



**HAL**  
open science

# A multi-level, multilingual approach to the annotation and representation of speech prosody

Daniel J. Hirst

► **To cite this version:**

Daniel J. Hirst. A multi-level, multilingual approach to the annotation and representation of speech prosody. Jonathan Barnes & Stefanie Shattuck-Hufnagel. Prosodic Theory and Practice, MIT Press, pp.117-149, 2022, 9780262543170. hal-03596616

**HAL Id: hal-03596616**

**<https://hal.science/hal-03596616>**

Submitted on 3 Mar 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## A multi-level, multilingual approach to the annotation and representation of speech prosody.

Daniel Hirst

### 1. Introduction: Prosodic annotation

Instinctively, when we listen to an utterance, we make a difference between *what* is being said, and *how* it is said, i.e. the prosody of the utterance. There is, generally, a fairly large consensus among linguists (and non-linguists) on the annotation of what is said. Interestingly, though, there is a remarkable lack of consensus on the annotation of how it is said. This difference can partly be explained by the fact that most writing systems provide a fairly reliable annotation of what is said, whereas there is very little written indication of how an utterance is spoken, apart from a limited number of abstract prosodic characteristics which are indicated by punctuation. The result is that there has been a bewildering variety of annotation systems for speech prosody - until quite recently, in fact, practically every researcher has developed their own annotation system. It was very rare for anyone to use someone else's system. Even for the 'British School', which had a fairly consistent approach to intonation patterns, there have been quite considerable differences between the annotation used by different authors such as Palmer (1924); Armstrong & Ward (1926); Jassem (1952); Kingdon (1958); O'Connor & Arnold (1961); Faure (1962); Halliday (1967); Crystal (1969); Couper-Kuhlen (1986) or most recently Wells (2006).

There have been two partial exceptions to this general tendency. One exception was the ToBI annotation system (Silverman et al., 1992), developed as a standard annotation system for American English, and later applied to several different languages (Jun, 2005). Another exception was my INTSINT system, described in more detail below. The original version of the system (Hirst, 1987) was based on an inventory of minimal pitch contrasts found in published descriptions of intonation patterns of sev-

---

<sup>1</sup>o

eral different languages. The aim was to provide a tool for the systematic description of these intonation patterns, something on the lines of the International Phonetic Alphabet. The INTSINT system was later used for nine of the twenty-two intonation systems described in Hirst & Di Cristo (1998).

## 2. Levels of representation

I have suggested (Hirst et al., 2000) that one of the reasons for this lack of consensus is that different annotation systems are in fact annotating different *levels* of representation. A similar point is made by Taylor (2009) who notes that different authors have completely different approaches, phonological, phonetic, acoustic:

These different approaches are of course quite common in all areas of language study; what is of note is that intonation seems to be one of the few areas where engineering and scientific approaches still have enough commonality to constitute a single field. [p 244]

The prosody of an utterance can be described at several different levels. In particular, I have argued that it is important to make a distinction between prosodic forms and prosodic functions. Prosodic forms are how the prosody affects what the utterance sounds like. Prosodic functions are what the prosody *does* and how it contributes to the interpretation of an utterance. I mentioned above that the only written correspondance to speech prosody is punctuation, which encodes a limited number of abstract prosodic characteristics. The difference between a full-stop, a comma, an ellipsis, a question mark and an exclamation mark, for example, corresponds to different ways in which the utterance is intended to be interpreted. As an example:

(1) a. OK.   b. OK,   c. OK...   d. OK?   e. OK!

Here, the first is intended as a statement, the second as an unfinished statement, the third as a tentative statement, the fourth as a question and the fifth as an emphatic statement. These are clearly examples of prosodic functions and the forms by which these functions are encoded may vary considerably from language to language and even from dialect to dialect.

The example of *OK* is particularly convenient for this, since the expression is used in numerous languages throughout the world, so the representations in (1) are not necessarily specific to a particular language.

This distinction between prosodic forms and prosodic functions<sup>2</sup> is, in my view, one of the sources of the disparity between prosodic annotation systems. In a ToBI annotation, such as, for example:

- (2) a. O K !  
b. % H\*L-L% //

the symbols H and L are clearly annotating prosodic form, while the symbols \*, -, and % seem much closer to indicating prosodic functions: whether or not a syllable is prominent, whether there is a boundary before what follows.

Interestingly, Wightman (2002) reported that there was considerably higher inter-transcriber agreement over the place of accents and boundaries than over the nature (H or L) of those accents and boundaries. In terms of our distinction between prosodic forms and prosodic functions, this can be interpreted as showing that listeners are much more proficient at identifying prosodic functions than prosodic forms. This seems a reasonable result, since it is well known that speakers and listeners are far more proficient at performing linguistic tasks (i.e. interpreting utterances) than meta-linguistic tasks like describing the forms of utterances.

It is also an interesting result in the light of the fact that for computers, the results are exactly the opposite: machines are excellent at describing the form of utterances but rather poor at interpreting them.

Prosodic forms, however, do not represent a single level of representation. In Hirst et al. (2000), I suggested that we need to distinguish at least three levels of prosodic form. Most abstract would be a level of underlying phonological representation, which would be the level at which the prosodic form of an utterance is linked to its prosodic function. Between this underlying phonological representation and the acoustic signal, I believe that we can usefully distinguish two more levels - a surface phonological representation and a phonetic representation.

---

<sup>2</sup>Faure (1962) makes a similar distinction between 'les caractères' (forms) and 'le rôle' (functions).

The term ‘phonetic representation’ has been used to cover a number of different phenomena which really need to be distinguished. Phonetics is sometimes used as a synonym for the acoustics and physiology of speech. It should be clear that I wish to distinguish these levels of analysis. Phonetics is also sometimes used as a synonym for ‘surface phonology’ as in the terms ‘phonetic transcription’ or the ‘International Phonetic Alphabet’.

My use of the term ‘phonetic’ follows that of Trubetzkoy (1949), for whom the distinction between phonology and phonetics is one between discrete and continuous phenomena. In this sense, then, a ‘phonetic transcription’ would more appropriately be termed a ‘surface phonological transcription’.

The level of surface phonology is a level of distinctive discrete categories with which we can describe surface phenomena cross-linguistically. The level of phonetics is the level of continuously variable phenomena from which we have factored out universal constraints on the production and perception of sounds. There has been much recent discussion about the ‘interface’ between phonetics and phonology. In the light of the above it would seem rather that phonetics is itself an interface – the interface between cognitive representations (phonology) and the physical manifestations of sound, perception and articulation.

Neither of these two levels of representation should necessarily be thought of as cognitive representations. It is quite possible that, in a complete theory of prosodic representation, neither of these levels would be necessary. Nevertheless, both levels appear to be useful as heuristic descriptive tools, particularly in the task of describing the prosody of languages which have not yet been adequately studied.

In the rest of this paper I first give a brief sketch of a functional prosodic annotation system which I developed many years ago for English intonation (Hirst, 1977). I then discuss the relationship of this annotation system to the notion of prosodic structure. After this, I work back up from the acoustic signal via the levels of phonetic representation and surface phonological representation, before giving some consideration to the nature of underlying phonological representation of prosody.

Most of the examples I give are in English and a few in French. These are two languages which are geographically close but prosodically rather different. There is nothing, however, in the framework which I describe which is specific to either of these two languages although of course it remains an empirical question whether the framework could be applied to

other languages without substantial modification.

### 3. Prosodic functions

In Hirst (1977), a revised and translated version of my unpublished PhD thesis (Hirst, 1974), I proposed a system of functional annotation to account for some of the ways in which intonation contributes to the meaning of an utterance. The book, written just a few years after Chomsky & Halle (1968) *Sound Pattern of English*, proposed a functional description of English intonation by means of a set of five distinctive features, for which I coined the term *intonative features*. These features were: [ $\pm$  STRESS;  $\pm$  CENTRE;  $\pm$  EMPHASIS;  $\pm$  BOUNDARY;  $\pm$  TERMINAL]. The features were transcribed using the following symbols:

- (3) a. 'no   b. °no   c. no   d. no /   e. no +   f. no //

The symbol in (3d.) was used at the beginning of an utterance or where the terminal/non-terminal nature of the boundary was not relevant, while (3e.) and (3f.) correspond respectively to [ $+$  BOUNDARY; -TERMINAL] and [ $+$  BOUNDARY; +TERMINAL].

These features were each justified by examining minimal pairs where only a difference of the value of the feature in question corresponded to a clear difference in meaning. In most cases the ambiguity is not complete and contexts can be imagined where the interpretations given here could be reversed, but it still seems that, at least for speakers of British English, the most obvious interpretation of each utterance corresponds to the one given here and that this is triggered by the feature described.

The feature [ $\pm$  STRESS]<sup>3</sup>, for example, was justified by the minimal pair:

- (4) a. he 'ate a little 'pudding (= he ate a small quantity of pudding)  
b. he 'ate a 'little 'pudding (= he ate a small pudding)

The feature [ $\pm$  CENTRE] (also known as *intonation centre*, *nuclear stress*, *primary sentence accent*, *tonic syllable* etc.) can be justified by the minimal pair:

---

<sup>3</sup>For reasons which I discuss below, I would, today, prefer to call this [ $\pm$  ACCENT] rather than [ $\pm$  STRESS].

- (5) a. it's a 'good 'job °too. (= It's a good thing too)  
 b. it's a 'good °job, °too (= The job is also a good one)

The feature [ $\pm$  EMPHASIS] can be justified by the pair:

- (6) a. if you °give it to me, I'll °mend it (= I'll mend it for you)  
 b. if you °give it to me, I'll °mend it (= I'll only mend it if you let me keep it)

The feature [ $\pm$  BOUNDARY] (i.e. *intonation unit boundary*) can be justified by the following pair:

- (7) a. she 'should have °phoned / her 'mother was °worried  
 (= her mother)  
 b. she 'should have °phoned her / 'mother was °worried  
 (= our mother)

Finally the feature [ $\pm$  TERMINAL] can be justified by the pair:

- (8) a. / 'would you like °tea + or °coffee + (= or something else?)  
 b. / 'would you like °tea + or °coffee// (= which?)

My conclusion, in the presentation of these 'features', and which I still believe today, was that:

the five intonative features we have described are sufficient to account for all the syntactic ambiguities we have so far come across and which are disambiguated by intonation, including a considerable number of ambiguities which up to now have always been treated from a semantic, "attitudinal" point of view"  
 [p 44]

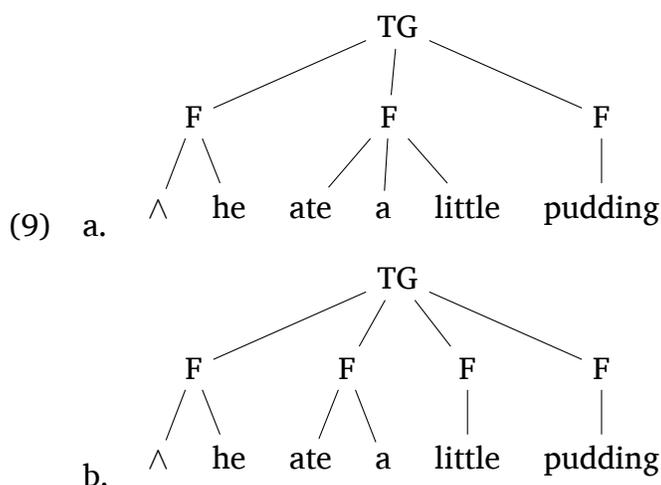
This does not, of course, mean that these five features exhaust the number of different intonational meanings which can be conveyed by prosody. They do, however, to my mind, constitute a minimal 'core' system of annotation for prosodic functions which any theory of intonational meaning should be able to account for.

#### 4. Prosodic structure

Since I wrote my thesis and book, a lot of research has been done, particularly in the framework of autosegmental-metrical phonology, on the nature of prosodic structure (Lieberman & Prince, 1977; Selkirk, 1984; Beckman & Pierrehumbert, 1986; Nespor & Vogel, 1986; Goldsmith, 1990; Ladd, 1996). It has today become a common point of view that both stress and boundaries are nothing but the reflection of a prosodic structure which is, at least to some extent, independent from syntactic structure. This seems a more convincing approach than a distinctive feature analysis.

##### 4.1. The Stress Foot

Halliday (1967), following Abercrombie (1964) had already suggested that an English utterance is a hierarchical phonological structure, containing four layers: the *Tone Group*, the *Foot*, the *Syllable* and the *Phoneme*. With this model, examples like (4) could be represented phonologically (omitting the level of the syllable and phoneme) as:

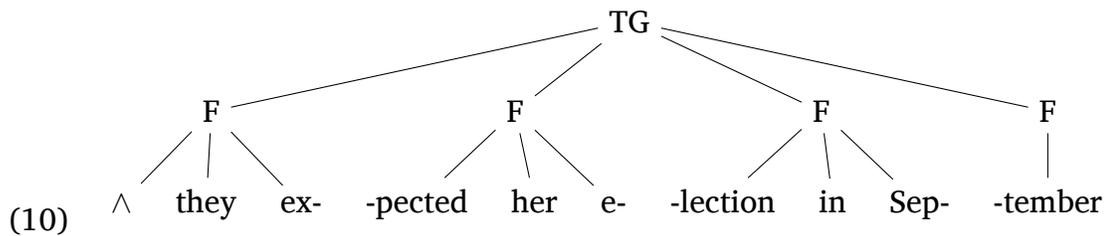


The symbol [ ^ ] in the first foot is what Halliday calls a “silent ictus”, corresponding to a rhythmic beat occurring before the first syllable of the foot, similar to the way that a *rest* is used in musical annotation for a similar purpose.

The Abercrombie/Halliday *Foot*, often called the *Stress Foot* to distinguish it from the more restricted phonological unit of autosegmental-metric phonology, can be defined for speech as a sequence of syllables

beginning with an accented syllable or with a silent beat at the beginning of a sentence, and continuing up to (but not including) the next accented syllable or the end of the *Tone Group*.

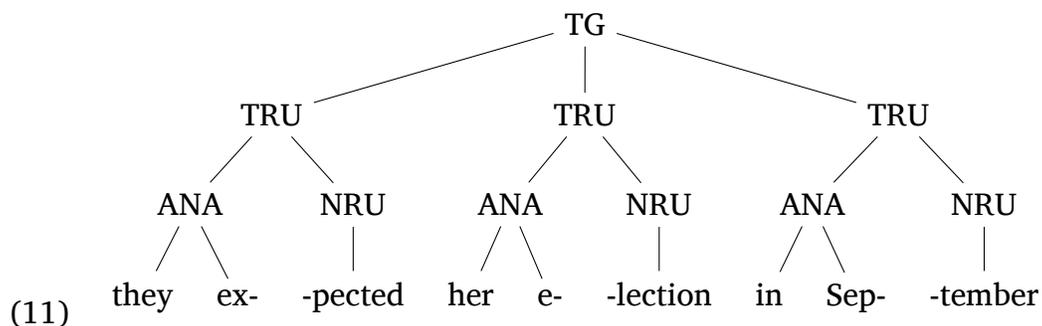
In (9), the foot boundaries all correspond to word boundaries, but this is not always necessarily the case. In (10), *none* of the foot boundaries corresponds to a word boundary, except of course the last.



#### 4.2. *Tonal Units and Rhythm Units*

Unlike Abercrombie and Halliday, who use the same unit to describe both melody and rhythm, Jassem (1952), more than ten years earlier, had made a clear distinction between what he called the *Tonal Unit*, which is conceived of as the domain of occurrence of local pitch movements in English, on the one hand, and the *Narrow Rhythm Unit* and the *Anacrusis*, on the other hand, conceived of as the domain of segmental timing and rhythm. The *Tonal Unit*, in Jassem’s model, was in fact identical to the unit which was thirteen years later called the *Foot* by Abercrombie as defined above.

The *Narrow Rhythm Unit* is similar to the *Foot*, except for the fact that it does not usually cross word boundaries, except in cases of enclitics, which are treated prosodically as belonging to the previous word. Any syllables which are not part of a *Narrow Rhythm Unit (NRU)*, form an *Anacrusis (ANA)*, in which, according to Jassem, the syllables are “pronounced extremely rapidly” (p 40). The *Anacrusis*, together with the following *Narrow Rhythm Unit*, together make up what Jassem termed the *Total Rhythm Unit (TRU)*. Example (10) in this analysis would look like (11):



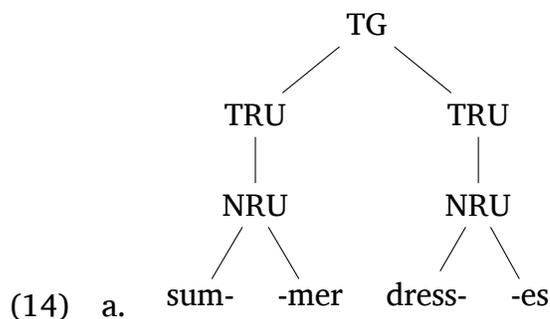
Jassem gave several examples where his representation allowed subtle distinctions of rhythm to be made which could not be captured without the notion of anacrusis. In the following minimal pair, for example:

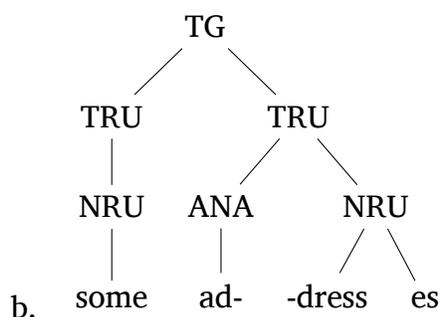
- (12) a. Summer dresses    b. Some addresses

Jassem notes that although the phonemes and stresses are identical, there is a subtle difference of rhythm between the two, the first syllable of (12a) being shorter than that of (12b) whereas the second syllable is longer in (12a) than in (12b). He attributes this difference to the fact that first two syllables of (12a) constitute a single NRU, whereas in (12b) the first syllable constitutes a NRU on its own and the second syllable constitutes an Anacrusis. He proposed to represent this in the phonetic transcription by the simple device of a space after each NRU as in:

- (13) a. /'sʌm ə 'dresɪz/    b. /'sʌm ə 'dresɪz/

corresponding to the tree representations:



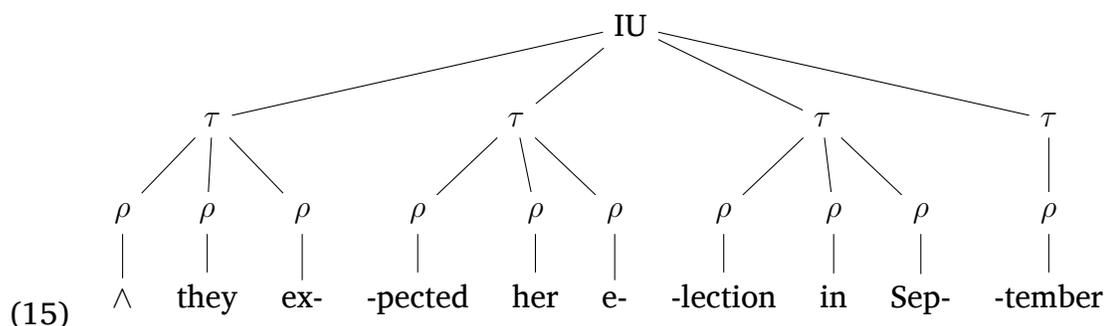


In a recent collection of articles, published to celebrate Wiktor Jassem’s 90th birthday, I suggested (Hirst, 2012) that Jassem’s *Total Rhythm Unit* does not actually play any phonological role. Instead we can combine the *Tonal Units* and the *Rhythm Units* (*Anacrusis* and *Narrow Rhythm Units*) into a single representation. In fact, if we do that, we can note that there is no longer any need to make a formal difference between *Anacrusis* and *Narrow Rhythm Unit*, since the *Narrow Rhythm Unit* will always coincide with the beginning of a *Tonal Unit* - both can simply be characterised as *Rhythm Units*. It is here that the distinction between stress and accent plays a crucial and, I believe, particularly interesting role.

Following Bolinger (1958), I interpret *accent* as referring to the physical manifestation of prominence in an utterance and *stress* as referring to the syllable or syllables of a word that are lexically marked and normally carry this prominence when the word is accented.

In my interpretation of Jassem’s model, then, the *Tonal Unit* (which I label [ $\tau$ ]) begins with an *accented* syllable, or at the beginning of the *Intonation Unit* (which I label [ $IU$ ], = Jassem’s *Tone Group*), and continues until the next accented syllable or until the end of the *Intonation Unit*. Deviating slightly from Jassem’s original formulation, but not, I hope, from the spirit of his analysis, I now define a *Rhythm Unit* (which I label [ $\rho$ ]) as beginning either at the beginning of a word, or with a *stressed* syllable and ending at the end of the word or before a stressed syllable. Once again as mentioned above, some cases of enclitics should be considered as forming part of the preceding word.

With this formulation, we can now combine the two levels of representation into a single prosodic tree as in (15):



Jassem's categories can, of course, be derived from this representation: a *NRU* is the first  $\rho$  in a  $\tau$ , excluding [ $\wedge$ ], and all other consecutive  $\rho$ 's can be grouped to form Jassem's anacrusis.

It can be seen that the rhythm unit  $\rho$  and the tonal unit  $\tau$ , thus defined, form an interface between the lexical representation of the word and the actual pronunciation of the utterance, or between stress and accent as defined by Bolinger (op. cit.).

In a study of the segmental duration of a large corpus of English, Hirst & Bouzon (2005) found that, as predicted by Jassem but contrary to Halliday's model, word boundaries *do* play an important role in the rhythmic structure of English. Strong negative correlations were found between the duration of a segment and the number of phonemes in the stress-foot, in the narrow rhythm unit and in the word, but no similar effect was found either in the syllable or in the anacrusis. Moreover the correlation was greater for the NRU than either the stress-foot or the word, thus confirming Jassem's predictions.

#### 4.3. *Emphasis and Terminality*

It seems, then that the prosodic functions corresponding to *Stress/Accent*, *Centre* and *Boundary* can more adequately be accounted for by an explicit model of prosodic structure such as I have just outlined.

The remaining two functions I mentioned: *Emphasis* and *Terminality*, do not look as if they can be accounted for by a difference in prosodic structure. One possibility which I have tentatively explored (Hirst, 1983, 1998) is that these features could be accounted for as prosodic morphemes, whose function is similar in many ways to sentence particles as found in many languages. Wakefield (2009, 2012) develops this idea more explicitly with a direct comparison between Cantonese sentence particles and English intonation patterns.

I will assume here, for the sake of argument, that English contains (at least) two abstract prosodic morphemes which I will represent as [E] for an emphatic particle and [Q] for a question particle, since one of the major functions of [ $\pm$ TERMINAL] is expressing prosodically a distinction between a statement and a question. If we want to distinguish between the intonation of questions and that of unfinished statements, a debatable issue as I discuss in (Hirst, 1998), we could suppose a third particle for that.

Question particles are well known in the description of many languages both tonal and non-tonal, e.g. *ma* in Mandarin Chinese and *li* in Russian. Emphatic particles are also common. In Bambara (Mali) (Bird et al., 1977), a sentence like:

- (16) a. Muso fila bé Sísé fe  
b. Sissé has two wives

can be modified by the emphatic particle *de* giving:

- (17) a. Muso de fila bé Sísé fe  
b. Sissé has two WIVES

- (18) a. Muso fila de bé Sísé fe  
b. Sissé has TWO wives

or

- (19) a. Muso fila bé Sísé de fe  
b. SISSE has two wives

where each time the emphatic particle *de* immediately follows the emphasized word.

In section (8), I return to the question of how these abstract particles may be interpreted prosodically.

## 5. Phonetic representation

### 5.1. Models of fundamental frequency

The search for an appropriate scale for measuring fundamental frequency has been one part of a systematic attempt, in particular by researchers from Holland ('t Hart et al., 1990), to develop a model of the way in which pitch is perceived. This was done by stylising raw fundamental frequency patterns as a sequence of straight lines, such that when the stylised frequency is used to resynthesise the utterance, the result is judged to be perceptually equivalent to the original intonation pattern.

Another approach to modelling pitch has been to attempt to model the way in which pitch is produced by speakers. In particular, work by Fujisaki and his colleagues has applied a model of pitch production (Fujisaki, 1991) to a large number of languages, including several tone-languages, analysing an intonation pattern as the superposition of three underlying components: a global baseline component, a sequence of phrasal components and a sequence of shorter accent components. These three components are added in the logarithmic domain to produce a raw fundamental frequency curve.

A third approach has been to develop acoustic models which are neither directly models of speech perception nor of speech production but which are compatible with both. This is the approach which I follow in this presentation.

### 5.2. Micro-melodic effects

Fitting a raw  $f_0$  curve with a mathematical model is not a simple straightforward problem due to the fact that fundamental frequency curves are not always continuous: unvoiced portions of the utterance have no associated  $f_0$ . Even when the curve is continuous it is often not smooth and this type of irregularity is often hard to model simply.

If we look at a two second extract of an utterance, corresponding to the words “More news about the Reverend Sun Myung M(oon)”<sup>4</sup>we notice (Figure 1(a)) that the beginning and end of the  $f_0$  curve is in fact fairly continuous and smooth. The reason for this becomes obvious when we look at the phonemes associated with the curve as in Figure 1(b).

---

<sup>4</sup>From the first recording (A01) of the Aix-Marsec corpus (auran-et-al2004aixmarsec)

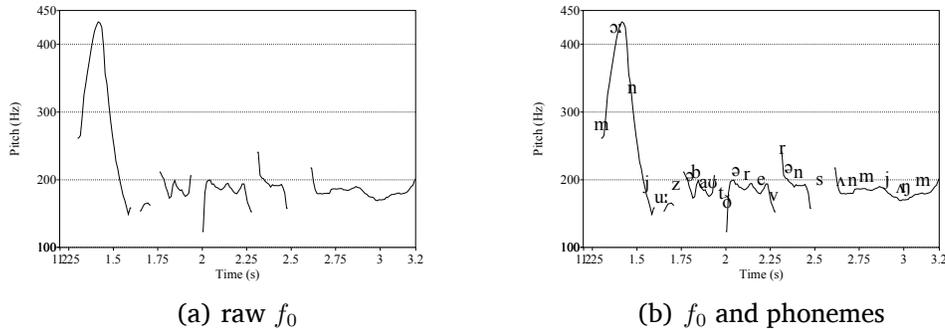


Figure 1: Two second extract of the  $f_0$  curve corresponding to “More news about the reverend Sun Myung M(oon)”.

The smooth portion at the beginning of the extract corresponds to the phonemes /mɔːnjuː/ whereas that at the end of the extract corresponds to /ʌnmjʌŋm/. All of these phonemes are fully sonorant, either vowels, semi-vowels or nasals. The discontinuity and irregularity of the  $f_0$  curve is due to the presence of obstruents in the utterance: stops and constrictives, which either interrupt the curve (for voiceless obstruents) or make it irregular (for voiced obstruents). The effect of these consonants has been called *micromelodic* as distinct from the *macromelodic* characteristics of larger pitch movements associated with accents and intonation patterns. Micromelodic effects, then, are caused by the aerodynamic characteristics of the articulation of different phonemes. Phonemes like vowels and sonorants, which hardly obstruct the airflow, have virtually no micromelodic effect whereas stops and constrictives disturb or interrupt the flow.

The raw fundamental frequency curve then can be thought of as the interaction between two components, a micromelodic component which is conditioned by the segmental nature of the individual speech sounds and a macromelodic component which corresponds to the underlying laryngeal gesture. This corresponds to the observation made by Nooteboom (1997) that we do not perceive the observable discontinuities of raw pitch-patterns unless they are longer than about 200 ms., as if human perception unconsciously bridges the silent gap by filling in the missing part of the pitch contour.

Linguists have been aware for a long time that fundamental frequency curves obtained from utterances containing only sonorants and vowels

are much better behaved than raw  $f_0$  curves obtained from unrestricted speech. It is for this reason that linguists have often constructed sentences consisting of mainly sonorants and vowels such as Eva Gårding’s “Madame Marianne Mallarmé har en mandolin från Madrid” (Madam Marianne Mallarmé has a mandolin from Madrid) for Swedish (Gårding, 1998), Yukihiro Fujisaki’s ‘Aoi aoinoewa yamanouenoieni aru.’ (The picture of the blue hollyhock is in a house on top of the hill.) (Fujisaki, 2004) or Annti Iivonen’s “Laina lainaa Lainalla lainen” (Laina lends Laina a loan) for Finnish (Iivonen, 1998). (cf Figure 2).

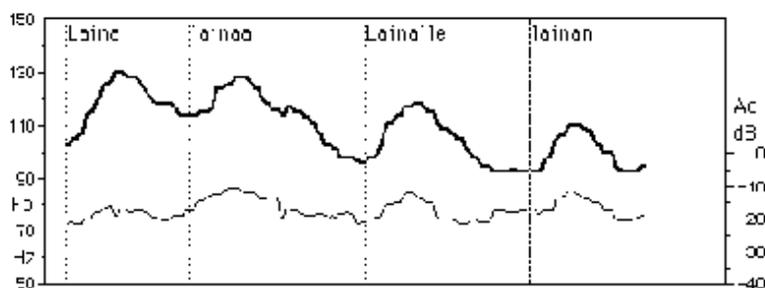


Figure 2: An example of a Finnish intonation pattern on a sentence with only sonorant phones. The sentence is “Laina lainaa Lainalla lainen” (‘Laina lends Laina a loan’) (Iivonen, 1998)

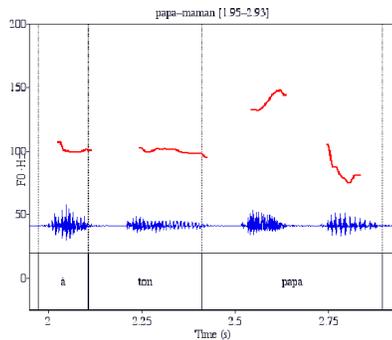
- . The top curve is the fundamental frequency and the lower curve is the intensity.

### 5.3. Macromelody and micromelody

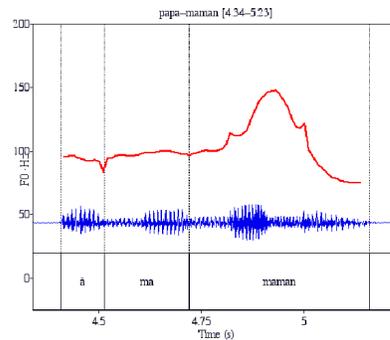
The idea, then, is that a raw intonation pattern is the interaction between two independent components: a macromelodic component determined by the accentuation and intonation of the utterance and a micromelodic component determined by the segmental phonemes.

If we compare two simple French phrases like “A ton papa.” (To your daddy.) and “A ma maman.” (To my mummy), pronounced as statements, we can see that there is the same underlying macromelodic pattern for the two utterances and that the surface differences are simply due to the different phonemes of the utterances, voiceless stops in Figure 3(a) and sonorant nasals in Figure 3(b).

What is particularly worth noting is that the  $f_0$  curve shown in Figure 3(a) is practically superposable on that of Figure 3(b). It seems as if the  $f_0$



(a) A ton papa. (To your daddy.)

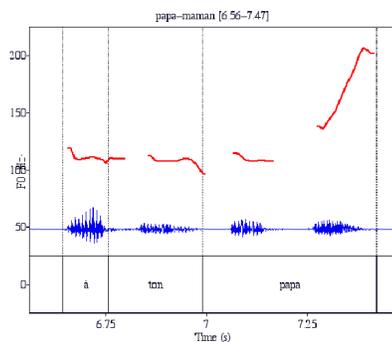


(b) A ma maman. (To my mummy.)

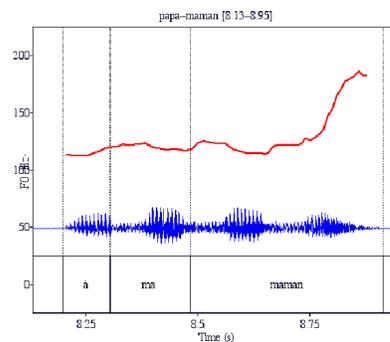
Figure 3: Two French sentences with a declarative intonation pattern

curve continues to change during the voiceless segments of the utterance even though it is not of course visible. This, of course, is not surprising if we think in terms of a continuous change of tension of the vocal folds, which can of course continue to change even during voiceless segments.

If we now compare these patterns with those observed on the same phrases pronounced as questions, as in Figure 4, we see once more that the two contours are practically superposable. And once again it seems clear that the  $f_0$  curve in Figure 4(a) continues to change during the voiceless segments of the utterance.



(a) A ton papa ? (To your daddy?)



(b) A ma maman ? (To my mummy?)

Figure 4: Two French sentences with an interrogative intonation pattern

In particular, the final rise on 'papa' does not begin at the onset of the final vowel: the  $f_0$  at this point is already considerably higher than that at

the end of the preceding vowel.

Notice that this idea of a continuously varying underlying pitch contour is not the model which is generally assumed in phonological descriptions of tonal and intonation contours. In the majority of these studies it is assumed that tones are directly associated with vowels, which are the ‘tone-bearing elements’ (cf Goldsmith (1990) p. 44 for example) and that the fundamental frequency observed on the consonants is simply an interpolation between the tones on the vowels. If that were so, then it might be thought that the pitch curve visible in Figures 3(a) and 4(a) is actually closer to the underlying form than that in Figures 3(b) and 4(b) which are simply the result of an interpolation on the sonorant consonants. The fact that the  $f_0$  curve follows the same trajectory in utterances with voiceless consonants as the smooth and continuous curve observed on the utterances with sonorants, however, and in particular the fact that the curve continues to evolve during the non-voiced portions of the utterance, seems to me convincing evidence that the planning of these curves is the result of an underlying macromelodic pattern on which the micromelodic variations are subsequently superimposed.

#### 5.4. *A model for $f_0$ curves*

##### 5.4.1. *Macromelodic and Micromelodic profiles*

The macromelodic component of an intonation pattern, then, has, I shall assume, the two characteristics of being smooth and continuous. This is fortunate because, as I mentioned above, modelling a discontinuous or irregular function is much more difficult than one which is continuous and smooth. To return to the example we saw earlier, the underlying macromelodic profile of the curve might be something like the continuous red curve in Figure 5(a), which could correspond to what we would produce if we were to hum the sentence instead of pronouncing it. I return below to how the continuous red curve was actually obtained; for the moment let us just assume that we have such a macromelodic profile.

Once we have a macromelodic profile, we can derive the micromelodic profile by dividing each value of the raw  $f_0$  curve by the corresponding value of the modelling function. This gives a result as displayed in Figure 5(b), where the smooth, continuous (red) curve is the macromelodic profile and the discontinuous (blue) curve is the micromelodic profile derived as just described.

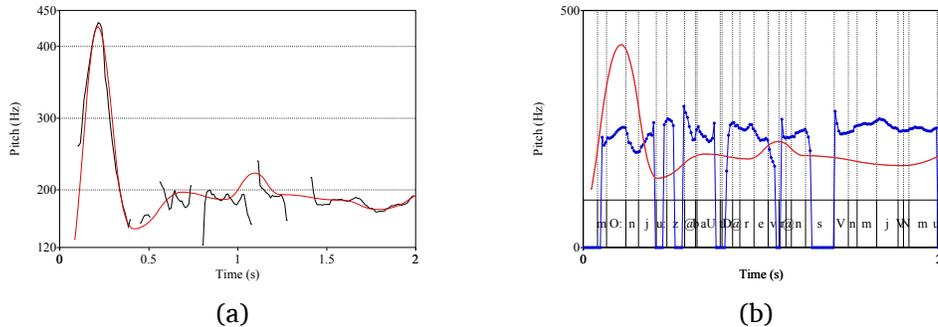


Figure 5: Raw  $f_0$  (black) together with macromelodic (red) and micromelodic (blue) profiles for the first two seconds of recording A01

Notice that such a modelling technique is not simply a stylisation of the raw  $f_0$  curve, since the raw curve has actually been factored into two orthogonal components without any loss of information. Multiplying each value of the red curve from Figure 5(b) by that of the blue curve in the same figure gives the original raw  $f_0$  as in the black curve in Figure 5(a). For speech synthesis, it would of course be possible to model the micromelodic profile itself and to use this to improve the segmental quality of the utterance (for an application to Arabic see Chentir et al. (2009)). For the study of intonation, the resynthesis of the utterance with the macromelodic profile is generally of sufficient high quality to evaluate the appropriateness of the modelled curve.

#### 5.4.2. $f_0$ transitions

One of the simplest ways to model a smooth continuous function like that in Figure 5(a) is as piecewise sequence of transitions between successive points on the curve. In previous work, until recently, I have referred to these points, following a fairly long tradition, as *target* points, but it should be noted that this name was not intended to imply that the ‘target points’ necessarily have any actual psychological reality for the speaker and listener. In order to make it clear that these points are not intended to represent cognitive targets, like those, for example in the model proposed in Xu & Wang (2001), I now prefer to use the term ‘anchor points’, which more clearly corresponds to their role in describing a pitch curve.

The advantage of a piecewise function over a global function is that each segment of the curve is defined locally by its own set of parameters,

which means that a modification of one portion of the curve does not entail modifications throughout the rest of the curve. The simplest model, of course, would simply be a linear transition between two anchor points as in Figure 6(a), where the transition is defined by the function

$$h_i = h_1 + \frac{(t_2 - t_i)}{(t_2 - t_1)} \cdot (h_2 - h_1) \quad (1)$$

where  $h_1$  and  $h_2$  are the  $f_0$  values of two adjacent anchor points and where  $t_1$  and  $t_2$  are the corresponding time values of these anchor points.

Here the  $f_0$  value ( $h_2$ ) of the second anchor point is higher than that ( $h_1$ ) of the first anchor point but the same reasoning would apply if it had been lower.

Naturally occurring  $f_0$  curves, of course, are not linear but curvilinear. A number of mathematical functions have been used in the past to model such functions. One of the simplest of these is a quadratic transition, corresponding to a constant acceleration followed by a constant deceleration of the pitch change as shown in Figure 6(b) and as defined by the functions:

$$t_i \in [t_1 \dots t_k] : h_i = h_1 + \frac{(h_2 - h_1) \cdot (t_i - t_1)^2}{(t_k - t_1)(t_2 - t_1)} \quad (2)$$

$$t_i \in [t_k \dots t_2] : h_i = h_2 + \frac{(h_1 - h_2) \cdot (t_i - t_2)^2}{(t_k - t_2)(t_1 - t_2)} \quad (3)$$

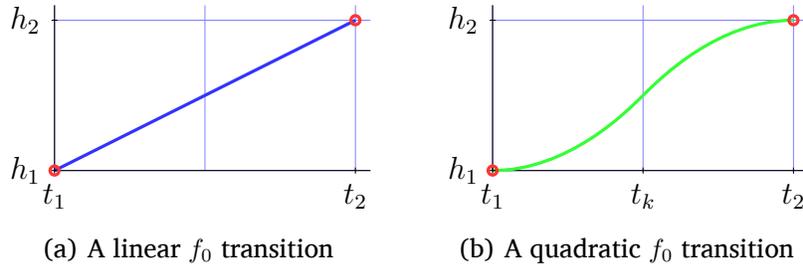


Figure 6: Linear and quadratic transitions from a first anchor point  $\langle t_1, h_1 \rangle$  to a second point  $\langle t_2, h_2 \rangle$ .

As can be seen from Figure (6(b)), in the case of a rise the transition consists of a concave curve from time  $t_1$  to time  $t_k$ , the point of maximum slope, followed by a convex curve from  $t_k$  to  $t_2$ .

Figure (7) shows the same extract we have seen several times now, with the anchor points which define the curve represented as green circles.

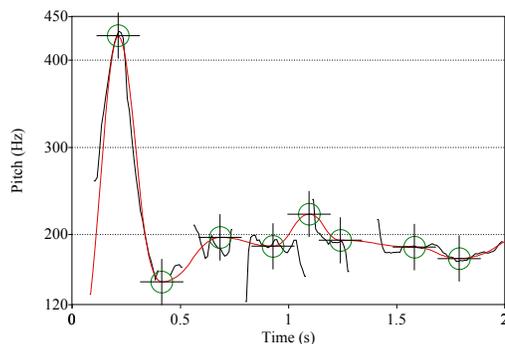


Figure 7: Macromelodic profile (red) for a two-second extract from recording A01, defined as quadratic transitions between anchor points (green).

### 5.5. *Momel*

A piecewise quadratic function such as that illustrated here is known as a quadratic *spline* function and has been in use in our laboratory since the 1980s to model intonation patterns using an algorithm called *Momel* (for “modelling melody”).

The *Momel* representation is, in fact, formally equivalent to a subset of the contours which can be produced by the Rise/Fall/Connection (RFC) model of intonation later developed by Paul Taylor (Taylor, 1994) as a tool for speech synthesis. The only difference is that the RFC model allows linear interpolations between two successive anchor points as well as quadratic interpolations. Of course if two successive anchor points have the same value of  $f_0$  then the transition will be linear (i.e. flat) with *Momel* too. It remains an empirical question whether there exist cases where a non-flat linear transition gives a better approximation to an  $f_0$  curve than a quadratic one. I have personally never observed a case.

The original implementation of *Momel* allowed the user to define anchor points manually by clicking on a representation of the  $f_0$  curve on the computer screen. The user could then resynthesise the utterance using PSOLA resynthesis. This can be done today with Praat (Boersma & Weenink, 1992 (2015)) by creating a Manipulation object and then removing and adding Pitch points manually. Praat displays the Pitch curve

with linear interpolation between the Pitch points but an approximation of the quadratic spline function can be obtained by the command **Interpolate quadratically...**, which was my own modest contribution to the Praat software.

Manual modelling of  $f_0$  is, of course, highly subjective and it was for this reason that my colleague Robert Espesser and I developed an automatic version of the algorithm (Hirst & Espesser, 1993), based on our experience of using the manual implementation of the model over a period of several years. The algorithm, which is described in detail in (Hirst et al., 2000) uses a form of robust regression to optimise the modeling of raw fundamental frequency curves with a quadratic spline function.

The algorithm was later evaluated on a corpus of read speech in 5 languages (corpus *Eurom1*) during the course of the Multext European project (Véronis et al., 1994). Evaluators were instructed to correct the anchor points for the modelled speech only when such corrections made an audible improvement to the resynthesis of the speech. Figure 1 shows the statistics for these corrections for the corpus of read speech together with those for a corpus of spontaneous speech in French.

Table 1: Results for the evaluation of the automatic Momel algorithm on read speech for 5 languages (English, French, German, Italian and Spanish and for spontaneous speech in French (corpus *Fref*). Columns show Total number of anchor points detected, number of anchor points added manually, number of anchor points deleted and the statistical measures of recall, precision and F-measure. Data from (Campione, 2001)

Corpus	Lang.	No. of points			Evaluation		
		<i>detected</i>	<i>added</i>	<i>deleted</i>	<i>recall</i>	<i>precision</i>	<i>F-measure</i>
<i>Eurom</i>	en	8380	623	125	93.0	98.5	95.7
	fr	6547	423	130	93.8	98.0	95.9
	ge	13595	1145	506	92.0	96.3	94.1
	it	9475	337	330	96.4	96.5	96.5
	sp	8985	651	16	93.2	99.8	96.4
<i>Fref</i>	fr	9835	532	744	94.5	92.4	93.4

The results of the evaluation were very encouraging. The F-measures for the different languages showed a global efficiency of around 95% and even the corpus of spontaneous French showed an F-measure of 93.4%

even though the algorithm had not at all been optimised for spontaneous speech.

Examination of the errors in the anchor points showed that one type of error in particular occurred systematically. This concerned a pitch rise before a silent pause where, frequently, the algorithm missed the final rise entirely. An example of this type of error can be seen in Figure (8). This error is understandable since the algorithm uses a local modelling technique, fitting a parabola to portions of the curve. As can be seen in the example, the raw  $f_0$  exhibits a concave portion corresponding to the acceleration of the pitch rise but hardly any convex portion corresponding to the decelerating portion of the rise. This is, in fact, the reason why the original algorithm fails to produce a final anchor point.

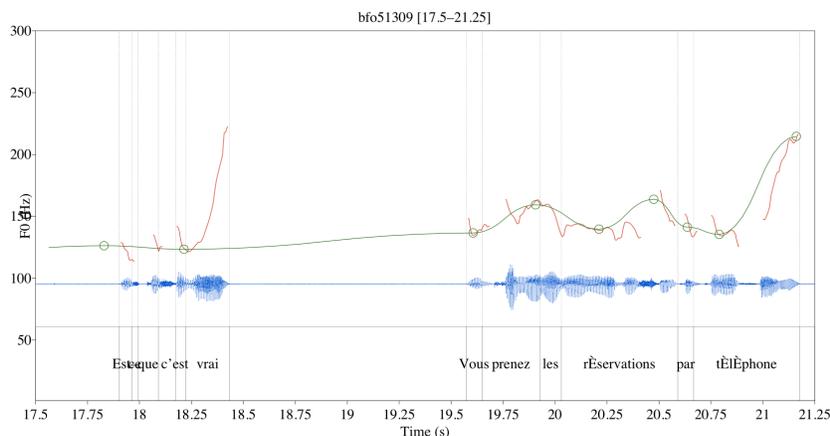


Figure 8: Old version of the automatic Momel algorithm for the utterance “Est-ce que c’est vrai ? Vous prenez les réservations par téléphone ?”. Raw (red) and modeled  $f_0$  (green).

The Momel algorithm has since been implemented as a Praat plugin Hirst (2007), which allows users to use its functions directly from the Praat menus without needing to handle scripts directly. The systematic error we have just seen has been corrected by a special treatment before silent pauses. The concave part of the pitch rise is now extended to include a similar shaped convex portion. In other words, in order to obtain the best fit for this pitch rise, a high anchor point is calculated, situated in the silent pause and as near as possible to the previous low anchor point. The anchor point is calculated so that the concave portion of the pitch

rise follows the raw data points as closely as possible. The result of this improved algorithm can be seen in Figure 9.

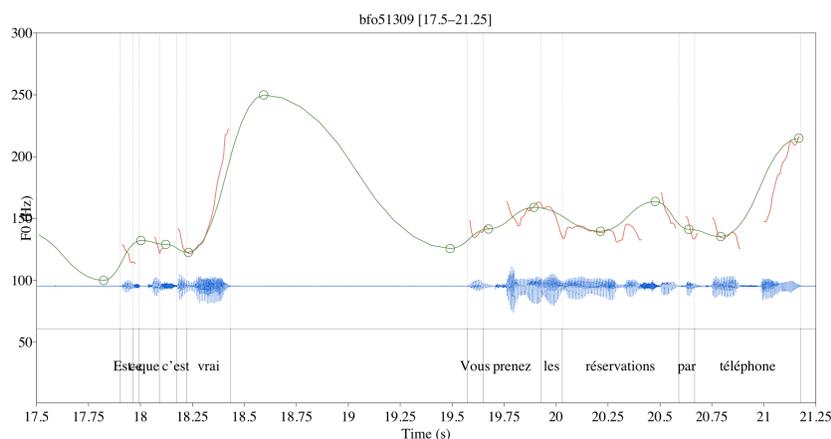


Figure 9: New version of the automatic Momel algorithm for the utterance “Est-ce que c’est vrai ? Vous prenez les réservations par téléphone ?” Raw (red) and modeled  $f_0$  (green).

An evaluation of the improved algorithm was carried out on a corpus of read speech in Korean (Hirst et al., 2007). It showed a significant and systematic improvement as compared to the older version of the algorithm. It is, naturally, desirable that the modelling tools we use should be as theory-neutral as possible. Complete neutrality, though, is obviously not entirely feasible, since any model necessarily makes some assumptions about the nature of underlying representations, as we saw above in the discussion of whether the underlying contour should be based only on the contours observed on the vowels or whether it should be modelled as a continuous underlying contour.

Rather than suggest that Momel is theory-neutral, then, I like to think that it is what we could call *theory-friendly*. I believe, that is, that the algorithm can be compatible with a number of different theoretical approaches to the description of speech melody. It has, in fact, been used in the past as a first step for modelling with the Fujisaki model (Mixdorff, 1999). It has also been used as first step for ToBI for both English (Maghbouleh, 1998; Wightman & Campbell, 1995) and Korean (K-ToBI) (Cho & Rauzy, 2008). It is also, of course, compatible with our own surface phonological representation alphabet, INTSINT which I describe below.

## 6. Surface phonological representation

### 6.1. *INTSINT: an International Transcription System for INTonation*

Momel, as we have seen, provides a reversible modelling of the raw  $f_0$  with no loss of information. INTSINT, an INternational Transcription System for INTonation, was originally developed as a tool for linguists to provide a surface phonological representation of an intonation pattern. The original version of the system (Hirst, 1987) was based on an inventory of minimal pitch contrasts found in published descriptions of the intonation patterns of numerous languages. The aim was to provide a tool for the systematic description of these intonation patterns, something along the lines of a narrow transcription using the International Phonetic Alphabet (IPA). Like the IPA, it was intended that INTSINT could be used for preliminary descriptions of intonation patterns, even for languages which had not previously been described. Notice that this aim is very different from that of the ToBI system (Silverman et al., 1992), which pre-supposes that the inventory of intonation patterns for the language being described has already been established. The official website for ToBI makes this particularly explicit:

Note: ToBI is not an International Phonetic Alphabet for prosody. Because intonation and prosodic organization differ from language to language, and often from dialect to dialect within a language, there are many different ToBI systems, each one specific to a language variety and the community of researchers working on that language variety.  
(source: <http://www.ling.ohio-state.edu/tobi/>)

The INTSINT system (whose name was suggested to us by Hans ‘t Hart in a personal communication) was presented in Hirst & Di Cristo (1998) and was used there for the annotation of the intonation of 9 different languages. Basically, it describes an intonation contour as a sequence of Tonal segments which are labelled using an alphabet of 8 symbols. The tonal segments are assumed to be of three types:

**Absolute tones** *t*(op) *m*(id) *b*(ottom): These are assumed to refer to the corresponding position of the speaker’s current pitch range.

**Relative tones** *h*(igher) *s*(ame) *l*(ower): Unlike absolute tones, relative tones are assumed to be refined with respect to the preceding tonal segment.

**Iterative relative tones** *u*(pstepped) *d*(ownstepped): These are also defined relative to the preceding tonal segment but generally involve smaller pitch changes and often occur in a sequence of steps either upwards or downwards.

In the chapters in Hirst & Di Cristo (1998), the INTSINT tones were represented by graphic symbols represented between two horizontal lines and aligned with the text. These symbols were: Top [ $\uparrow$ ]; Bottom [ $\downarrow$ ]; Higher [ $\uparrow$ ]; Same [ $\rightarrow$ ]; Lower [ $\downarrow$ ]; Upstepped [ $<$ ] and Downstepped [ $>$ ]. The Mid tone was reserved for the unmarked onset of an Intonation Unit and was not marked.

In most later publications the capital letters T, M, B, H, S, L, U, D were used instead of the graphic symbols. Since Hirst (2011), I prefer to represent the INTSINT tones with lower case letters rather than upper case. This may help to avoid confusion with other more abstract coding schemes such as ToBI (Silverman et al., 1992), or the even more abstract underlying representation used in (Hirst, 1998) and section (8) below, both of which use some of the same letters as INTSINT.

## 6.2. Mapping from INTSINT to Momel

Although INTSINT was introduced as a descriptive tool for linguists, since its introduction was later than the creation of the  $f_0$  modelling tool, I already had in mind the possibility that this *surface phonological* annotation could be linked to the analysis of the  $f_0$  curve as a sequence of *phonetic* anchor points. I anticipated, then, that it might be possible to map the output of the *Momel* algorithm onto a sequence of symbols from the *INTSINT* alphabet. To do this, following the idea of an analysis by synthesis paradigm, it was first necessary to define a mapping in the other direction, that of synthesis. Some of the history of the way in which this mapping was defined is described in Hirst (2005).

In its current implementation, the mapping depends on two speaker/utterance-specific parameters called *key* and *span* which together define the speaker's

pitch range<sup>5</sup>. The key (like a musical key) defines a central reference point for the speaker's pitch range and the span defines the maximum and minimum pitch values of the range which are taken to be symmetrical (on a logarithmic scale) above and below the speaker's key.

The two parameters together define three absolute tones - top, mid and nottom - with respect to the speaker's pitch range as in the following formulas which assumes that the value of key is given in Hertz and the value of span in octaves:

$$t = key \times \sqrt{2^{span}} \quad m = key \quad b = \frac{key}{\sqrt{2^{span}}} \quad (4)$$

The pitch anchor points corresponding to the relative tones are then defined with respect to both the preceding anchor point (here called  $p$ ) and the top ( $t$ ) or bottom ( $b$ ) of the range.

An anchor point coded  $h$  is simply defined as the geometric mean (i.e. the mean on a log scale) of the preceding anchor point and the top of the range - it thus corresponds to a pitch movement which moves up halfway towards the value of  $t$ . As can be expected, an anchor point coded  $s$  is defined as the same as the preceding anchor point. Symmetrically to  $h$ , an anchor point coded  $l$  is defined as the (geometric) mean of the preceding anchor point and the bottom of the range.

$$h = \sqrt{p * t}; \quad s = p; \quad l = \sqrt{p * b} \quad (5)$$

For the iterative tonal segments  $u$  and  $d$  the implementation defines the values as the (geometric) mean of the value of the preceding anchor point and that which would be obtained if the anchor point were coded  $h$  or  $l$  respectively. In other words these anchor points correspond to a pitch excursion one quarter of the way to the top/bottom of the pitch range.

$$u = \sqrt{p * \sqrt{p * t}}; \quad l = \sqrt{p * \sqrt{p * b}} \quad (6)$$

---

<sup>5</sup>In some earlier publications I used the word *range* for what I here call *span*. I now prefer to use 'span' to refer to the interval, independently from the value of the 'key'. The two values 'key' and 'span' together define the speaker's 'range'. So we might say that a given speaker has a pitch range from 100 Hz to 200 Hz, corresponding to a span of one octave and a key of 141 Hz.

Assuming, once again, a logarithmic scale for the pitch range, these values are illustrated graphically in Figure 10.

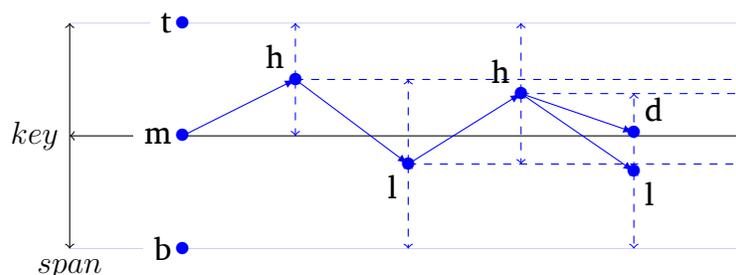


Figure 10: Graphic illustration of the mapping from INTSINT to Momel defined by 2 parameters *key* and *span*

This implementation obviously makes a number of assumptions, most of which would be open to empirical investigation. This is one of the major advantages of an explicit model such as this.

One consequence of the model, which was not specifically intended but which turns out to be fortunate, is that a sequence of alternating *h* and *l* tones will automatically introduce an iterative lowering of the tones, much like the probably universal effect of *downdrift* that has been described in the literature on tone and intonation as occurring in languages throughout the world. This downdrift effect, as illustrated in Figure (11), is thus an automatic by-product of the way in which the relative tones are defined.

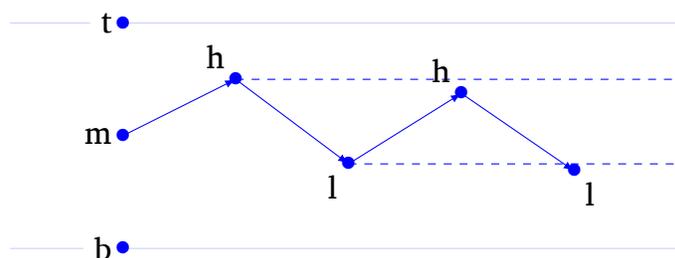


Figure 11: A graphic illustration of the fact that downdrift is an automatic by-product of the way in which the relative tones are defined

Since it was first developed, the Momel algorithm has been applied relatively successfully to a number of different languages, including English,

French, Italian, Catalan, Brazilian Portuguese, Venezuelan Spanish, Russian, Arabic, isiZulu and Korean (for references see Hirst (2007)). More recently (Zhi et al., 2010), the algorithm was applied to a corpus of speech in Standard (Beijing) Chinese. This was particularly challenging, since the corpus used was spontaneous speech and involved a language with a rich lexical tone system. An attempt to optimise window-size for the algorithm showed no overall improvement with respect to the manually corrected data. This was taken to confirm the fact (as had been suggested by Xu & Sun (2002)) that pitch change in a lexical tone language like Chinese is not notably faster than in languages with no lexical tone. The annotated data obtained during this application will constitute a useful yardstick for evaluating improvements to the automatic algorithm, which is expected to be far more robust than data annotated for languages with no lexical tone.

### 6.3. Mapping from Momel to INTSINT

Having defined a mapping from tonal segments (INTSINT) to pitch anchor points (Momel), the same model can be used to establish a reverse mapping from the anchor points to the tonal segments. As is generally the case, such a reverse mapping from continuous variables to discrete categories is much less straightforward than the mapping from categories to continuous variables. The approach we have adopted is an exhaustive search of the target space for the optimal values of the two parameters *key* and *span* together with the optimal coding of the sequence of anchor points, given those two parameters.

The procedure as described in Hirst (2005) has been implemented as a Perl script. It assumes that the relevant target space is defined as follows:

$$key = mean \pm 50Hz \quad (step : 1) \quad span = 0.5 \dots 2.5 \text{ octaves} \quad (step : 0.1) \quad (7)$$

The script thus tries each of the possible values of the two parameters within this target space. For each of the 2000 possible couples  $\langle key, span \rangle$  the script evaluates every possible coding of the anchor points using the formulas in equations (4, 5, 6) and calculating the sum of the square of differences between the predicted value and the observed value. The output of the script is thus the optimised value (within the target range) of the parameters *key* and *span* together with the optimal INTSINT coding

using these parameters. The output of the script is a text file such as the following corresponding to the application of the script to the same extract which we saw earlier:

```
; A01_01.intsint created on Tue Aug 24 08:12:47 2010 by intsint.pl 2.11
; from A01_01.momel
;   32 values mean = 191
<parameter span=1.4>
<parameter key=235>
0.113 M 221 235
0.219 D 205 208
0.434 D 182 190
0.746 B 120 145
1.177 S 120 145
1.423 T 428 382
1.623 B 146 145
1.894 U 197 184
...
```

Figure 12: Extract of a sample output of the automatic INTSINT coding script. After the values for *key* and *span* each line gives the time value (in seconds), the optimised INTSINT code, the original anchor point used as input and the predicted anchor point derived from the coding with the current values of *key* and *span*

In this output, the optimised values of *key* and *span* (here 235 Hz and 1.4 octaves respectively) are given, together with the sequence of anchor points. Each line gives the time value (in seconds) for the anchor point, the optimised INTSINT code, the original anchor point used as input and the predicted anchor point derived from the coding with the current values of *key* and *span*

Unlike with Momel, there is some loss of information with the INTSINT coding (as can be seen from the differences between the 3rd and 4th columns of the output in Figure (12) the model is still however a reversible one which can be used for synthesis. Figure (13) shows the result of resynthesising the pitch anchor points from the results of the automatic INTSINT analysis applied to a 5 sentence passage from the Eurom1 corpus, compared to the anchor points obtained from the application of the Momel algorithm.

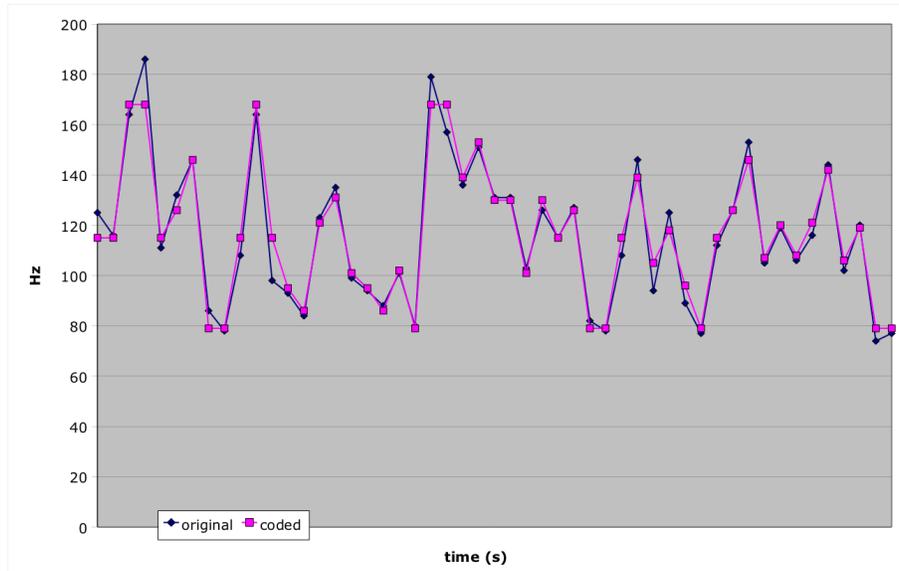


Figure 13: Predicted (pink) versus observed (blue) values for pitch anchor points for a 5 sentence passage after automatic coding with the INTSINT alphabet

#### 6.4. Longer term characteristics of pitch range

The implementation of INTSINT as described above presupposes that there are no variations in *key* and *span* within the segment of speech which is analysed. In authentic speech, such changes naturally occur quite frequently and are obviously very significant and important. One solution is to implement the INTSINT coding on smaller segments of speech such as breath groups, making what seems a fairly reasonable assumption that changes of *key* and/or *span* are more likely to occur *between* breath groups rather than *within* them. For an investigation of the possibility of automatising such a process, cf. De Looze (2010) which implements an algorithm (AdoReVa) applying a cluster analysis for this task.

### 7. ProZed

The availability of explicit models of speech melody such as those described above makes it possible to use such models to implement more abstract representations of prosody. With this aim in mind, my current

work in progress concerns the implementation of a *Prosody Editor for Linguists* called *ProZed* (Hirst, 2015) and which aims to provide linguists with a tool allowing them to experiment with different abstract phonological models, providing them with an acoustic output with which they can, at least informally, evaluate the relative value of different models.

A version of this editor applied to speech rhythm is described in Hirst & Auran (2005), which implemented an empirical linear model for rhythm where each *rhythm unit* is characterised by three parameters:  $q$  a (possibly speaker dependent) parameter defining a unit of quantal lengthening,  $t$  a long term parameter of *tempo*, and  $k$  a local scalar effect of lengthening specific to each rhythm unit. With these parameters the duration of the rhythm unit  $\rho$  is defined as:

$$\hat{d}_\rho = t \cdot \left\{ \sum_{i=1}^m \bar{d}_{i/p} + k \cdot q \right\} \quad (8)$$

where  $\bar{d}_{i/p}$  is the mean duration of all the phones labelled as the same phoneme  $p$  as that occurring in position  $i$  in the rhythm unit.

ProZed has been implemented as a plugin to the Praat software. It allows the manipulation of the rhythmic and the tonal aspects of speech as defined on three specific tiers in addition to the *phoneme* tier. The first two of these are named the *RU* (rhythm unit) tier and the *TU* (tonal unit) tier. These two tiers control the short term variability of segmental timing and tonal variation, respectively. Longer term variations of rhythm and melody can be controlled via a third tier named the *IU* (intonation unit) tier.

The speech input to the program may be natural recorded speech, the prosodic characteristics of which will then be modified by the software, or, alternatively it may be the output of a speech synthesis system with, for example, fixed durations for each speech segment.

The current version of the program is designed as the re-synthesis step of what is planned to be a complete analysis by synthesis cycle. This will be directly integrated with the output of the Momel-INTSINT and ProZed Rhythm analysis models which are described above and interfaced with the automatic alignment system SPPAS (Bigi & Hirst, 2012 (submitted, 2013)).

## 8. Underlying phonological representation

We have seen that it is possible to derive automatically from the acoustic signal a phonetic representation of an intonation contour (Momel), and that this can then be converted automatically into a time-aligned sequence of discrete symbols (INTSINT) which together with the output of the automatic alignment software can constitute at least a first approximation to a surface phonological representation of the prosody of an utterance.

It now remains to be seen how we can bridge the gap between the abstract functional representation of prosody which I outlined in section (3) and this surface phonological representation.

For British English, the classical description of the intonation of non-emphatic declarative utterances and WH-questions, as in Armstrong & Ward (1926); Kingdon (1958); O'Connor & Arnold (1961) is of a pitch pattern which begins on a mid-level with a rise to a high pitch on the first accented syllable, followed by a step-like lowering of pitch on each accented syllable until the last accented syllable followed by a fall to the bottom of the speaker's pitch-range.

Among the examples of this 'tune', Armstrong and Ward (op.cit.) give:

(20) It was the 'last 'thing I ex- 'pected to 'find there.

In terms of prosodic structure as described above in section (4), this utterance can be parsed into one intonation unit containing six tonal units:

(21) [ [^ it was the ] [last] [thing I ex-] [-pected to] [find there] ]

The INTSINT annotation of the tonal representation of this utterance (ignoring here the question of the relative alignment of the tones and the segmental material) could then be:

(22) [m It was the [h last ] [d thing I ex- ] [d -pected to ] [d find there  
b] b]

where the initial [m] and final [b] tones are attached directly to the intonation unit and the other tones are attached to the tonal units.

Pike (1945) describes a neutral intonation pattern for declarative utterances in American English as a sequence of falling contours associated with each accented syllable. This could be annotated as:

(23) [m it was the [h last l] [h thing I ex- l] [h -pected to l] [h find there  
b] b]

Pike characterises a “descending stress series” such as in (22) as expressing “EXTREME precision, or certainty” (p 70). Because of this difference in interpretation, to American ears the British intonation may sound over-precise or pompous, while to British ears the American pattern may sound over-enthusiastic or exuberant.

This dialectal difference between American and British English intonation patterns is interesting since a downstepped tone in many African tone languages can often be traced to a historical low tone which is sometimes described as a “floating” tone (Clements & Ford, 1979), i.e. a tone which is not phonetically realised as a pitch target but which has the effect of lowering the following high tone.

The two dialects could consequently be described as having the same underlying representation with the difference that in British English the low tone is somehow “delinked” from the prosodic structure.

If we compare this intonation to that of the French translation of 21:

(24) C’était la 'dernière 'chose que je m’atten-'dais à 'trou ver 'là.

We are likely to find that instead of a sequence of falling pitch-movements, like in American English (as described by Pike (op. cit.) or a sequence of “downsteps” as in British English, the French utterance is produced as a sequence of rising pitch movements, culminating with the accented syllable and with a final falling pitch on the last syllable, cf. (Hirst & Di Cristo, 1984; Di Cristo, 1998).

We have suggested (Hirst & Di Cristo, 1984) that one of the major prosodic differences between English and French is that while, in English, Tonal Units consist of an accented syllable followed by a number of unaccented syllables, in French they consist of an accented syllable preceded by a number of unaccented syllables. The example (24) could consequently be annotated<sup>6</sup>:

---

<sup>6</sup>French words are regularly accented on the final syllable but non-final words may also be accented on the first syllable, especially when immediately followed by an accented syllable as in ‘dernière chose’ and ‘trouver là’ in this example. See Hirst & Di Cristo (op.cit.) and Di Cristo (op.cit.) for details.

(25) [[c'était la der-] [-nière chose] [que je m'attendais] [à trou-] [ver là]].

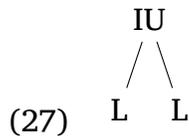
with which we can associate an INTSINT representation:

(26) [**m** [**m** c'était la der- **h**] [**l** -nière chose **h**] [**l** que je m'attendais **h**] [**l** à trou-] **h**] [**d** -ver là]**b**]

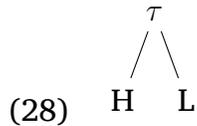
For an underlying phonological representation of these tunes, I will assume, following Pierrehumbert (1980) that there are only two values for underlying tones: L and H, and that the specific INTSINT interpretations of these underlying tones are allophonic variants determined by the context.

Under this interpretation, it seems clear, then, that for all three varieties we have looked at, American English, British English and French, the underlying tones associated directly with the Intonation Unit (= boundary tones) are [L L].

We can assume, then, a template like the following to account for this association:

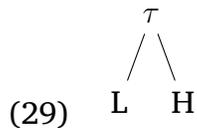


For the Tonal Units, we can assume that both British English patterns and American English patterns are derived from a single template:



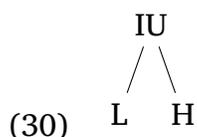
and that in British English, but not in American English, there is a downstepping rule that will have the effect of interpreting an underlying sequence [H L] [H ...] as [h] [d ...]

French patterns can be derived from:



assuming simply that, in French, a downstepping rule applies only to the last tonal unit, interpreting [... H] [L H] as [... h] [d].

For the intonation of yes-no questions in English and French, the only change which needs to be made to the underlying representation for statements is that the final L is replaced by H, triggered by the presence of the question particle [Q] so that we have:



For French, the presence of a question particle triggers not only the application of the template in (30) but also the application of a downstepping rule on each accent so that a sequence [... H] [L H] is interpreted as [... h] [d] for each Tonal Unit in the Intonation Unit.

For English, the presence of an emphatic particle [E] has two effects. The pitch accents on the accented syllables are either deleted entirely or considerably reduced, as if there were an intermediate prosodic constituent between the Tonal Unit and the Intonation Unit:

- (31) a. [ [John has] [two] [wives] ]  
           (= *John has two wives*)  
       b. [ [ [John has] [two] [wives] ] ]  
           (= *JOHN has two wives*)  
       c. [ [John has] [ [two] [wives] ] ]  
           (= *John has TWO wives*)  
       d. [ [John has] [two] [ [wives] ] ]  
           (= *John has two WIVES*)

There might, then, be a condition that only the first tonal unit within this emphatic constituent is assigned tones. This cannot be the whole story, though, since the last example of (31) would then be assigned the same tonal representation as the first, unemphatic example. In fact the pitch accent on an emphatic constituent is systematically higher than other

accents. One possibility, then, would be to assume that the emphatic constituent is assigned [H L] tones just like the tonal units, so the underlying representation of (31) would then be something like:

- (32) a. [L [H John has L] [H two L] [H wives L] L]  
 (= *John has two wives*)
- b. [L [H [H John has L] [two] [wives] L] L]  
 (= *JOHN has two wives*)
- c. [L [H John has L] [H [H two L] [wives] L] L]  
 (= *John has TWO wives*)
- d. [L [H John has L] [H two L] [H [H wives L] L] L]  
 (= *John has two WIVES*)

which could then be interpreted as the INTSINT annotations (for British English), where the **h** or **d** followed immediately by **h** would ensure a higher pitch on the following syllable than when there is just a single **h**:

- (33) a. [**m** [**h** John has] [**d** two] [**d** wives **b**] **b**]  
 (= *John has two wives*)
- b. [**m** [**h** [**h** John has **b**] [two] [wives] **b**] **b**]  
 (= *JOHN has two wives*)
- c. [**m** [**h** John has] [**d** [**h** two **b**] [wives] **b**] **b**]  
 (= *John has TWO wives*)
- d. [**m** [**h** John has] [**d** two] [**d** [**h** wives **b**] **b**] **b**]  
 (= *John has two WIVES*)

## 9. Conclusion

I have presented in this paper a sketch of a fairly complete model of prosody, although a great number of details concerning the implementation of the model need to be specified. The model extends from a functional annotation of intonation at the most abstract level, via an underlying

phonological representation and a surface phonological representation to a phonetic representation which is directly convertible into an acoustic signal. The model is described in more detail in my forthcoming book (Hirst, forthcoming), together with more detailed justification of the approach presented here. It is obvious that at each level I have made a number of arbitrary choices which may not stand up to further investigation, but I am nonetheless convinced that a multilevel, multilingual model of this type is what is needed to further our knowledge of the complex way in which prosody contributes to our interpretation of utterances.

## References

- Abercrombie, D. (1964). Syllable quantity and enclitics in English. In D. Abercrombie, D. Fry, P. MacCarthy, N. Scott, & J. Trim (Eds.), *In Honour of Daniel Jones* (pp. 216–222). London: Longman.
- Armstrong, L., & Ward, I. (1926). *A Handbook of English Intonation*. Cambridge: Heffer.
- Beckman, M., & Pierrehumbert, J. (1986). Intonational structure in Japanese and English. *Phonology yearbook*, (pp. 255–309).
- Bigi, B., & Hirst, D. (2013). What's new in SPPAS 1.5? In *Proceedings of Tools and Resources for the Analysis of Speech Prosody* (pp. 62–65.). Aix-en-Provence, France.
- Bigi, B., & Hirst, D. J. (2012 (submitted)). SPEech Phonetization Alignment and Syllabification (SPPAS): a tool for the automatic analysis of speech prosody. In *Proceedings of the 6th International Conference on Speech Prosody*.
- Bird, C. S., Hutchinson, & Kante (1977). *Beginning Bambara (An Ka Ba-manankan Kalan)*. Indiana University Press.
- Boersma, P., & Weenink, D. (1992 (2015)). Praat, a system for doing phonetics by computer. <http://www.praat.org> [version 5.4.06, February 2015].
- Bolinger, D. (1958). A theory of pitch accent in English. *Word*, 14, 119–149.

- Campione, E. (2001). *Etiquetage semi-automatique de la prosodie dans les corpus oraux - algorithmes et méthodologies.*. Ph.D. thesis Université de Provence.
- Chentir, A., Guerti, M., & Hirst, D. (2009). Extraction of standard arabic micromelody. *Journal of Computer Science*, 5, 86–89.
- Cho, H., & Rauzy, S. (2008). Phonetic pitch movements of accentual phrases in korean read speech. In *Proceedings of the 4th International Conference on Speech Prosody*. Campinas Brasil.
- Chomsky, N., & Halle, M. (1968). *The Sound Pattern of English*. Coll. Studies in Language. New York, NY, USA : Harper & Row, 1968, 470 p.
- Clements, G., & Ford, K. (1979). Kikuyu Tone Shift and its Synchronic Consequences. *Linguistic Inquiry*, 10, 179–210.
- Couper-Kuhlen, E. (1986). *An introduction to English prosody.*. London: Edward Arnold.
- Crystal, D. (1969). *Prosodic systems and intonation in English*. Cambridge: Cambridge University Press.
- De Looze, C. (2010). *Analyse et interprétation de l'empan temporel des variations prosodiques en français et en Anglais*. Ph.D. thesis Université de Provence Aix-en-Provence, France.
- Di Cristo, A. (1998). Intonation in French. In D. J. Hirst, & A. Di Cristo (Eds.), *Intonaton Systems. A Survey of Twenty Languages*. (pp. 195–218). Cambridge University Press.
- Faure, G. (1962). *Recherches sur les caractères et le rôle des éléments musicaux dans la prononciation anglaise.*. Paris: Didier.
- Fujisaki, H. (1991). Modeling the generation process of F0 contours as manifestation of linguistic and paralinguistic information. In *Proceedings of the XIIth International Congress of Phonetic Sciences* (pp. 1–10).
- Fujisaki, H. (2004). Information, prosody and modeling - with emphasis on tonal features of speech. In *Proceedings of the 2nd International Conference on Speech Prosody* (pp. 1–10).

- Gårding, E. (1998). Intonation in Swedish. In D. Hirst, & A. Di Cristo (Eds.), *Intonation Systems. A Survey of Twenty Languages*. chapter 6. (pp. 117–136). Cambridge: Cambridge University Press.
- Goldsmith, J. A. (1990). *Autosegmental and metrical phonology*. Cambridge, Mass.: B. Blackwell.
- Halliday, M. (1967). *Intonation and Grammar in British English*. Mouton, 62p.
- 't Hart, J., Collier, R., & Cohen, A. (1990). *A Perceptual Study of Intonation: an Experimental-Phonetic Approach to Speech Melody..* Cambridge University Press.
- Hirst, D. J. (1974). *La levée de l'ambiguïté syntaxique par les traits intonatifs : essai de formalisation..* Ph.D. thesis Université de Provence.
- Hirst, D. J. (1977). *Intonative features: a syntactic approach to English intonation*. Mouton, 135p.
- Hirst, D. J. (1983). Interpreting intonation : a modular approach. *Journal of Semantics*, 2, 171–181.
- Hirst, D. J. (1987). La représentation linguistique des systèmes prosodiques : une approche cognitive. Thèse de Doctorat d'Etat (Habilitation Thesis), Université de Provence.
- Hirst, D. J. (1998). Intonation in British English. In D. J. Hirst, & A. Di Cristo (Eds.), *Intonation Systems. A Survey of Twenty Languages*. chapter 3. (pp. 56–77). Cambridge: Cambridge University Press.
- Hirst, D. J. (2005). Form and function in the representation of speech prosody. *Speech Communication*, 46, 334–347.
- Hirst, D. J. (2007). A Praat plugin for Momel and INTSINT with improved algorithms for modelling and coding intonation. In *Proceedings of the XVIth International Conference of Phonetic Sciences* (pp. 1233–1236). Saarbrücken.
- Hirst, D. J. (2012). Empirical models of tone, rhythm and intonation for the analysis of speech prosody. In D. Gibbon, D. J. Hirst, & N. Campbell

- (Eds.), *Rhythm, Melody and Harmony in Speech. Studies in Honour of Wiktor Jassem*. (pp. 23–33). Poznan: Polish Phonetic Association volume 14/15 of *Speech and Language Technology*.
- Hirst, D. J. (2015). ProZed: A speech prosody editor for linguists, using analysis-by-synthesis. In K. Hirose, & J. Tao (Eds.), *Speech Prosody in Speech Synthesis: Modeling and generation of prosody for high quality and flexible speech synthesis*. chapter 1. (pp. 3–17). Berlin Heidelberg: Springer Verlag.
- Hirst, D. J. (forthcoming). *Speech Prosody: from Acoustics to Interpretation*. Springer Verlag.
- Hirst, D. J., & Auran, C. (2005). Analysis by synthesis of speech prosody: the prozed environment. In *Proceedings of Interspeech 2005. (Lisbon)* (pp. 3225–3228).
- Hirst, D. J., & Bouzon, C. (2005). The effect of stress and boundaries on segmental duration in a corpus of authentic speech (British English). In *Proceedings of Interspeech/Eurospeech 05*. (pp. 29–32.). Lisbon.
- Hirst, D. J., Cho, H., Kim, S., & Yu, H. (2007). Evaluating two versions of the momel pitch modeling algorithm on a corpus of read speech in korean. In *Proceedings of Interspeech* (pp. 1649–1652). Antwerp, Belgium volume VIII.
- Hirst, D. J., & Di Cristo, A. (1984). French intonation: a parametric approach. *Die Neueren Sprachen*, 83, 554–569.
- Hirst, D. J., & Di Cristo, A. (1998). *Intonation Systems: A Survey of Twenty Languages*. Cambridge University Press, 487 p.
- Hirst, D. J., Di Cristo, A., & Espesser, R. (2000). Levels of representation and levels of analysis for the description of intonation systems. In M. Horne (Ed.), *Prosody: Theory and Experiment. Studies Presented to Gösta Bruce*. (pp. 51–87). Kluwer Academic Pub.
- Hirst, D. J., & Espesser, R. (1993). Automatic modelling of fundamental frequency using a quadratic spline function. *Travaux de l'Institut de Phonétique d'Aix*, 15, 75–85.

- Iivonen, A. (1998). Intonation in Finnish. In D. Hirst, & A. Di Cristo (Eds.), *Intonation Systems. A Survey of Twenty Languages* chapter 17. (pp. 331–347). Cambridge University Press.
- Jassem, W. (1952). *Intonation of Conversational English: (educated Southern British)*. Nakl. Wroclawskiego Tow. Naukowego; skl. gl.: Dom Ksiażki.
- Jun, S.-A. (Ed.) (2005). *Prosodic Typology. The Phonology of Intonation and Phrasing*. Oxford University Press, London.
- Kingdon, R. (1958). *The groundwork of English intonation*. Longmans.
- Ladd, D. (1996). *Intonational Phonology*. Cambridge.: Cambridge University Press,.
- Lieberman, M., & Prince, A. (1977). On stress and linguistic rhythm. *Linguistic Inquiry*, 8, 249–336.
- Maghbouleh, A. (1998). Tobi accent type recognition. In *Proceedings of ICSLP*. Paper 0632.
- Mixdorff, H.-J. (1999). A novel approach to the fully automated extraction of fujisaki model parameters. In *Proceedings of ICASSP 1999*.
- Nespor, M., & Vogel, I. (1986). *Prosodic Phonology*. Dordrecht: Foris Publications.
- Nooteboom, S. (1997). The prosody of speech: melody and rhythm. In W. Hardcastle, & J. Laver (Eds.), *The Handbook of Phonetic Sciences*. (pp. 640–673). Oxford: Blackwell.
- O'Connor, J., & Arnold, G. (1961). *Intonation of Colloquial English. A Practical Handbook*. London: Longmans.
- Palmer, H. E. (1924). *English Intonation with Systematic Exercises*. Cambridge: W. Heffer and sons.
- Pierrehumbert, J. B. (1980). *The Phonology and Phonetics of English Intonation*. Ph.D. thesis Massachusetts Institute of Technology Cambridge, Mass.

- Pike, K. (1945). *The Intonation of American English*. Ann Arbor, Michigan: University of Michigan Press.
- Selkirk, E. (1984). *Phonology and Syntax: The Relation between Sound and Structure*. Current Studies in Linguistics I., 476 p. M.I.T. Press.
- Silverman, K., Beckman, M., Pitrelli, J., Ostendorf, M., Wightman, C., Price, P., Pierrehumbert, J., & Hirschberg, J. (1992). TOBI: A Standard for Labeling English Prosody. In *Second International Conference on Spoken Language Processing* (pp. 867–870). Banff, Canada.: ISCA.
- Taylor, P. (1994). The rise/fall/connection model of intonation. *Speech Communication*, 15, 169–186.
- Taylor, P. (2009). *Text-to-speech synthesis*. Cambridge University Press.
- Trubetzkoy (1949). *Grundzüge der Phonologie*. (French translation by J. Cantineau 1957) *Principes de phonologie*. Paris: Klincksieck.
- Véronis, J., Hirst, D., & Ide, N. (1994). NL and speech in the MULTEXT project. In *Proceedings of AAAI Workshop on Integration of Natural Language and Speech* (pp. 72–78). Seattle, USA.
- Wakefield, J. (2009). *The English Equivalents of Cantonese Sentence-final Particles: A Contrastive Analysis*. Ph.D. thesis Hong Kong Polytechnic University Hong Kong.
- Wakefield, J. (2012). A floating tone discourse morpheme: The english equivalent of cantonese lo1. *Lingua*, 122, 1739–1762.
- Wells, J. C. (2006). *English Intonation: An Introduction*. Cambridge: Cambridge University Press.
- Wightman, C. (2002). ToBI or not ToBI? In *Proceedings of the First International Conference on Speech Prosody* (pp. 21–29). Aix-en-Provence.
- Wightman, C., & Campbell, N. (1995). Improved labeling of prosodic structure. In *IEEE Transactions on Speech and Audio Processing*.
- Xu, Y., & Sun, X. (2002). Maximum speed of pitch change and how it may relate to speech. *Journal of the Acoustical Society of America*, 111, 1399–1413.

- Xu, Y., & Wang, Q. E. (2001). Pitch targets and their realization. evidence from Mandarin Chinese. *Speech Communication*, 33, 165–181.
- Zhi, N., Hirst, D., & Bertinetto, P. M. (2010). Automatic analysis of the intonation of a tone language. applying the momel algorithm to spontaneous standard chinese (beijing). In *Proceedings of Interspeech XI*. Makuhari, Japan.