

COMPONENTS OF PROSODIC EFFECTS IN SPEECH RECOGNITION

ANNE CUTLER

MRC Applied Psychology Unit
15 Chaucer Rd.
Cambridge CB2 2FF, U.K.

ABSTRACT

Previous research has shown that listeners use the prosodic structure of utterances in a predictive fashion in sentence comprehension, to direct attention to accented words. Acoustically identical words spliced into sentence contexts are responded to differently if the prosodic structure of the context is varied; when the preceding prosody indicates that the word will be accented, responses are faster than when the preceding prosody is inconsistent with accent occurring on that word. In the present series of experiments speech hybridisation techniques were first used to interchange the timing patterns within pairs of prosodic variants of utterances, independently of the pitch and intensity contours. The time-adjusted utterances could then serve as a basis for the orthogonal manipulation of the three prosodic dimensions of pitch, intensity and rhythm. The overall pattern of results showed that when listeners use prosody to predict accent location, they do not simply rely on a single prosodic dimension, but exploit the interaction between pitch, intensity and rhythm.

Speakers place accent on the most important words in an utterance. Thus by finding accented words, listeners can efficiently locate the most central parts of a speaker's message. Previous studies have shown that listeners do indeed actively use sentence prosody to tell them where accented words are going to occur. Cutler [2] produced pairs of sentences varying in prosodic contour. An example is (1):

- (1) (a) The couple had quarrelled over
a BOOK they had read.
(b) The couple had quarrelled over
a book they hadn't even READ.

Upper case represents sentence accent. In (1a) the main sentence accent falls on *book*, in (1b) on *read*. These sentences were used as materials in a phoneme-

monitoring experiment, in which listeners are asked to respond as quickly as possible to the presence of a specified word-initial phoneme. In (1), the target phoneme is /b/, so the target-bearing word is *book*. Targets on accented words are responded to faster than targets on unaccented words in this task. In Cutler's experiment, the target-bearing word itself was actually spliced out of both sentence contexts and replaced in each by identical copies of a neutral rendition of the same word. The result of this manipulation was a pair of sentences with acoustically identical target-bearing words, which were preceded by identical sequences of words; the only difference between the members of each pair was the prosody applied to the words preceding the target. In one case the prosodic contour in which the target-bearing word occurred was consistent with accent falling upon that word; in the other, it was consistent with the target-bearing word being unaccented. Under these conditions, the 'accented' targets still elicited faster responses than the 'unaccented' targets, and since the only relevant differences between the two sentences in each pair lay in the prosody, Cutler concluded that listeners must have used cues in the prosody to direct their attention to the location where sentence accent would fall.

Prosody, however, is not a unitary phenomenon. The separate dimensions of rhythm, pitch and intensity all contribute to the prosodic structure of an utterance. Cutler's experiment did not examine *how* listeners were exploiting prosody to predict accent, or whether any one prosodic dimension was more informative than others.

Cutler and Darwin [3] subsequently found that removing pitch information - i.e. monotonising the sentences - did not remove the accent effect; in monotonised spliced sentences like (1) the 'accented' targets are still responded to significantly faster than the 'unaccented' targets.

From this, Cutler and Darwin concluded that pitch information could not be a necessary component of the accent prediction effect. They speculated that no prosodic dimension might prove necessary for listeners to predict upcoming accents, but variation in any prosodic dimension might prove sufficient.

In the present studies, the three prosodic dimensions of pitch, rhythm and intensity are separately manipulated in an attempt to analyse the accent effect in further detail. Unlike the study by Cutler and Darwin, which simply removed the dimension of pitch by setting it to a single value across each utterance, the present studies investigate the effects of the separate prosodic dimensions when they are *interchanged* between the two members of a sentence pair. To begin with, using dynamic time-warping techniques in a system developed by Jeffrey Bloom at the Polytechnic of Central London [1], we exchanged the rhythmic patterns within each pair of sentences (for examples like [1], where naturally different contours were produced by having a slight variation in the text at the end of the sentences, the rhythmic patterns were exchanged up to the point at which the two members of the pair diverged). Thus (1a), for example, was given the rhythm of (1b) but retained its original pitch and intensity contours; (1b) had the rhythm of (1a) but its own pitch and intensity patterns.

In Experiment 1, phoneme-monitoring response times were measured in these rhythmically manipulated sentences, and in the same sentences with intact prosody. The intact sentences were LPC-analysed and resynthesised to control for acoustic effects of resynthesis. The words bearing the target were acoustically identical in all four sentences belonging to a set such as (1).

There were 20 such sentence sets. Forty listeners, in four groups of ten, took part in the experiment. Each group heard only one sentence from each set, and the two variables of 'accented' versus 'unaccented' targets, and intact versus rhythmically manipulated prosody, were counterbalanced across subject groups.

Subjects were tested individually. Response times, measured from a click (inaudible to the subjects) aligned with target onset, were collected by a microcomputer using programs developed by Norris [4]. After the experiment subjects were given a short recognition test, and their response times were analysed only if they scored at least two-thirds correct on this test.

The results of this experiment are shown in Fig. 1. The intact sentences, in which rhythm, pitch and intensity contours are preserved from the original utterance, show the advantage of 'accented' over 'unaccented' targets which was found in the earlier experiments. This indicates that the resynthesis alone is not interfering with listeners' ability to use prosodic contours to predict the location of accented. The difference in this condition is significant ($F(1,3.6) = 21.36, p < .001$). In the rhythmically manipulated sentences, however, the advantage of originally accented over originally unaccented targets is less than half as large as the difference in the prosodically intact sentences, and is not statistically significant ($F(1,3.6) = 3.55, p > .05$).

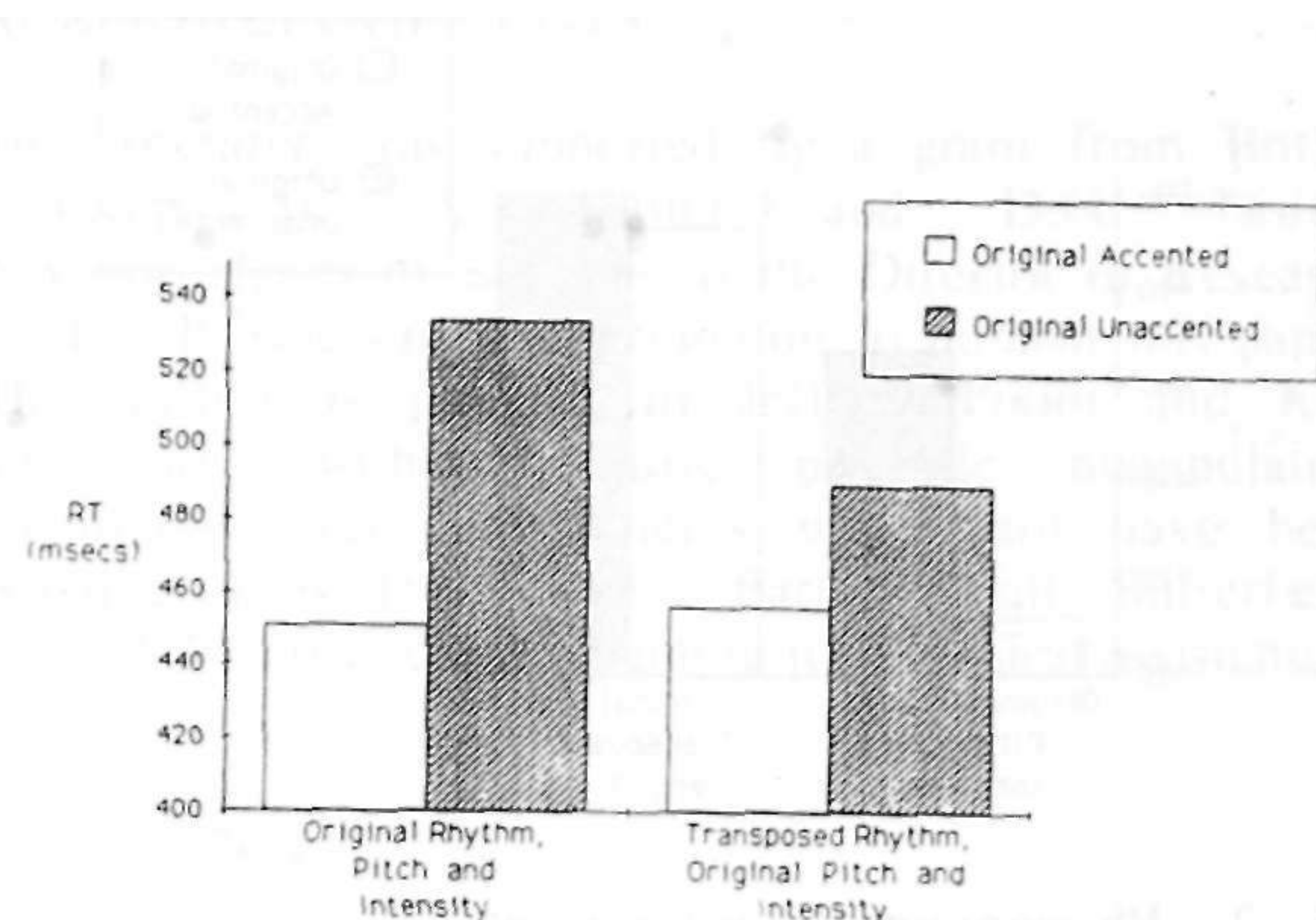


FIG. 1. Phoneme-monitoring response time (msecs.), Experiment 1.

This experiment shows that the rhythmic manipulations have severely affected the accent effect. Each of the utterances which had undergone this rhythmic manipulation had an unnatural, indeed a conflicting, prosodic structure - pitch and intensity contours signalled one prosodic pattern while the rhythm signalled another. It is clear that listeners did not base their prosodic processing on one aspect of the prosodic contour alone.

One possible interpretation of this result is that listeners are simultaneously processing all three prosodic dimensions, and that the separate contributions of each prosodic dimension to the predicted accent effect are simply additive. The attenuated, but still positive, effect in the rhythmically manipulated sentences would, on this simple story, be attributable to the combination of positive effects contributed by the pitch and intensity contours, set against a negative effect contributed by the rhythmic contour.

This interpretation was tested in Experiment 2. This experiment investigated prosodic manipulations which were the reverse of those in Experiment 1. The pitch and intensity contours were transposed between originally accented-target and unaccented-target members of a sentence pair, leaving the rhythmic contour, alone, intact.

This manipulation was possible because the time-warping applied to the sentences in Experiment 1 produced pitch and intensity contours which, although they preserved the contour shape from the utterance they had originally belonged to, were aligned with the rhythmic pattern of that utterance's pair. Therefore these contours could simply be transposed onto that pair. These transpositions were realised using prosodic editing routines devised by Kim Silverman.

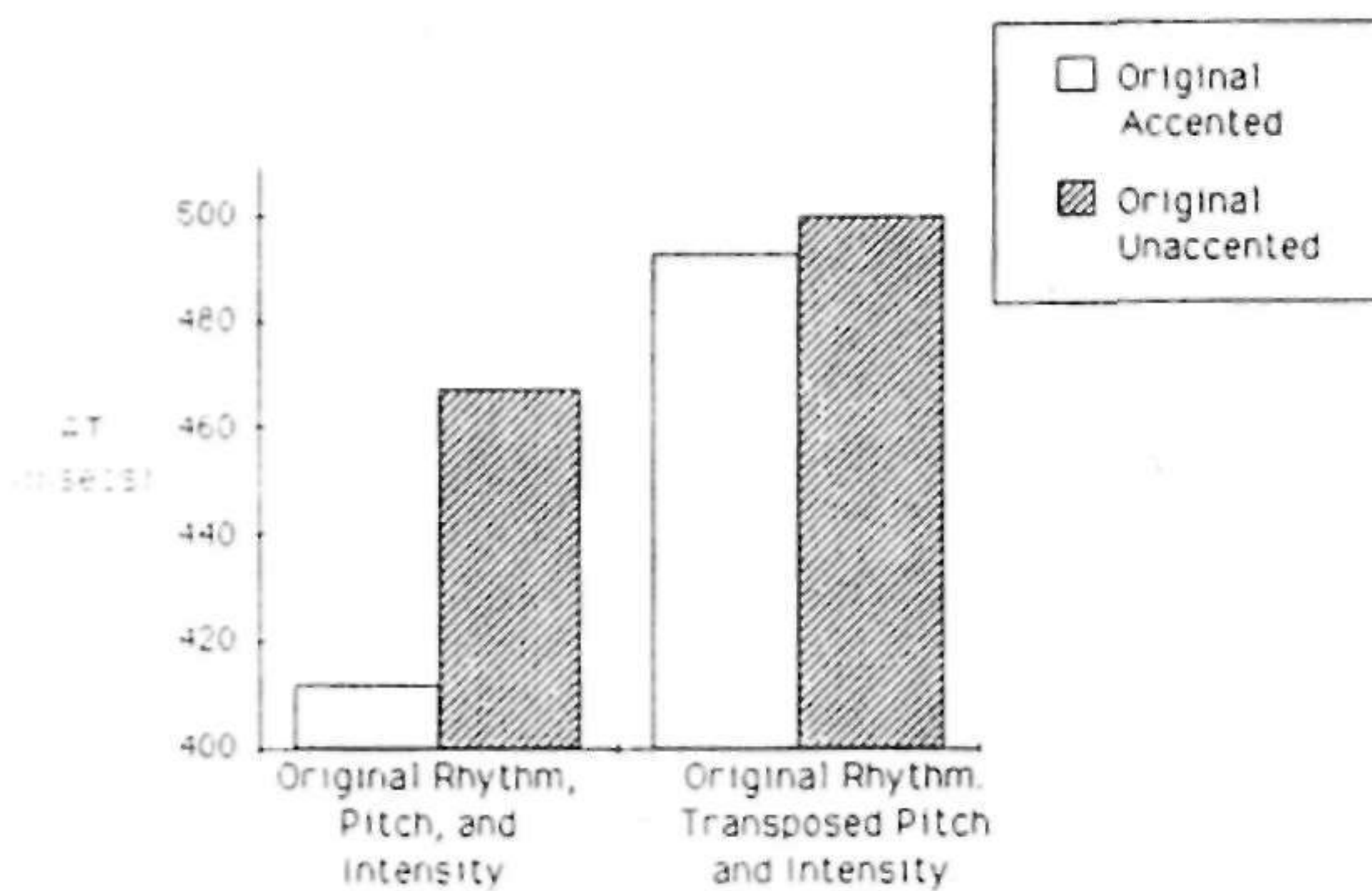


FIG. 2. Phoneme-monitoring response time (msecs.), Experiment 2.

Experiment 2, like Experiment 1, included the resynthesised utterances with intact prosody; these were compared with the utterances in which of the original prosody only the rhythm was preserved intact, the pitch and intensity contours being transposed between members of a pair. Again, the target-bearing words were acoustically identical in all sentences from any set.

Forty listeners, who had not taken part in Experiment 1, were tested; design and procedure were as in Experiment 1. The results are shown in Fig. 2.

It can be seen that once again the utterances with intact prosody showed a strong accent effect, i.e. response time advantage for 'accented' over 'unaccented' targets. This difference was statistically significant ($F(1,36) = 6.85$, $p < .02$). In the utterances with transposed pitch and intensity contours, there was virtually no response time difference between originally accented and originally unaccented targets ($F < 1$).

The results of this experiment rule out the very simple explanation of Experiment 1 offered above. Had listeners been simply evaluating all three dimensions of prosody in an additive fashion, we might have expected the reverse of the result found in Experiment 1 - that is, we might have expected an advantage of originally unaccented targets over originally accented targets of about half the magnitude of the difference in the opposite direction produced by the prosodically intact utterances. However, the conflicting prosody in this case wiped out any difference in response times as a function of original accent location.

This result raises the possibility that transposition of prosodic contours might itself interfere with listeners' ability to predict accent location by extracting relevant

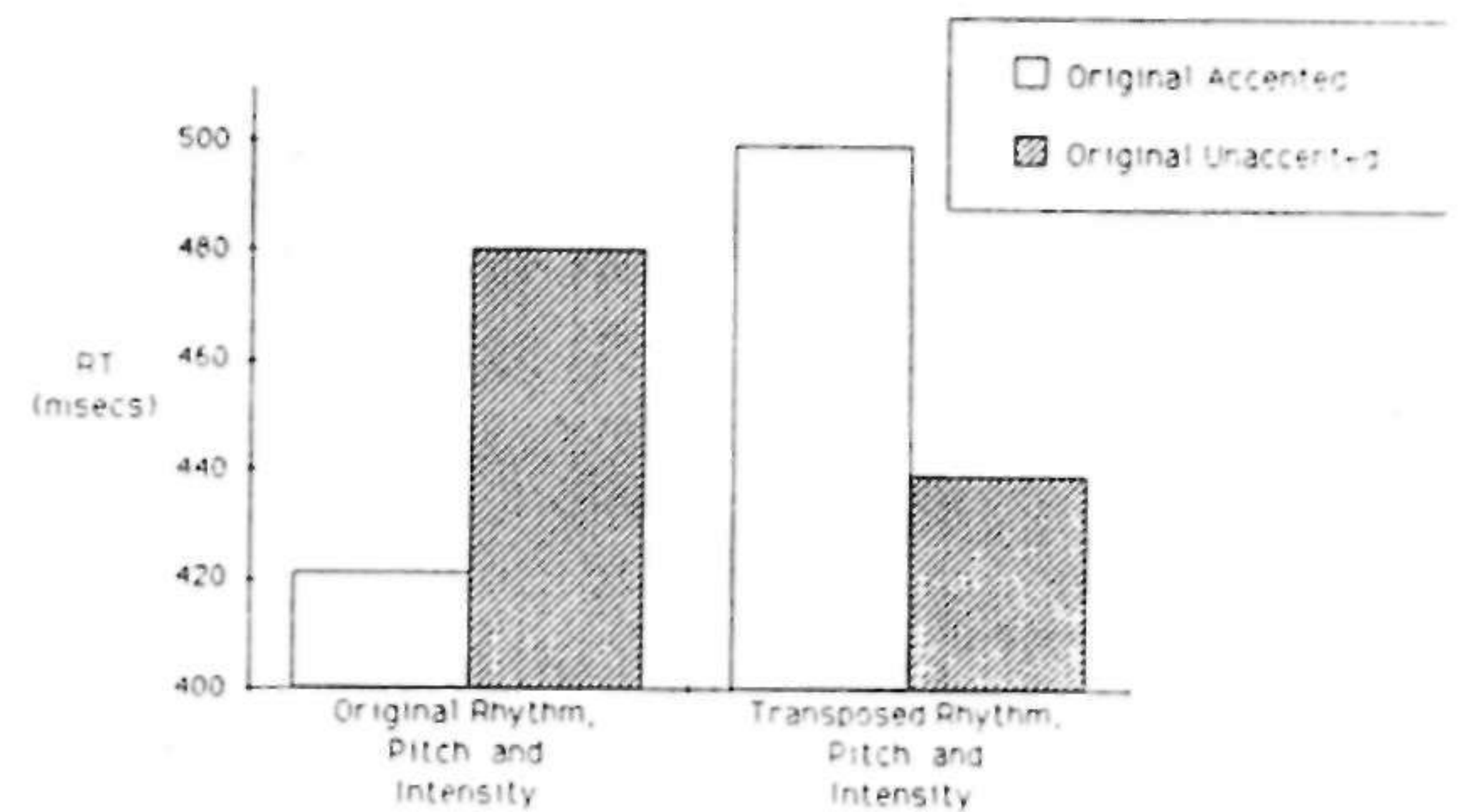


FIG. 3. Phoneme-monitoring response time (msecs.), Experiment 3.

information from the prosody. In order to rule out this possibility, a further experiment was conducted in which all three prosodic dimensions were transposed.

In Experiment 3, the resynthesised utterances with intact prosody were again tested, and compared in this case with utterances in which rhythm, pitch and intensity contours had all been transposed between members of a sentence pair. The manipulated utterances in this experiment therefore exhibited the maximum of transposition, in that every utterance had rhythm, pitch and intensity contours which had originally been applied to another utterance. However, they exhibited the minimum of prosodic conflict, since rhythm, pitch and intensity contours were always in accord.

As in the previous experiments, the target-bearing words were acoustically identical in all sentences from any set.

Forty listeners, none of whom had taken part in Experiments 1 and 2, were tested. Design and procedure were as in the preceding experiments. The results are shown in Fig. 3.

Once again there was a significant advantage for 'accented' over 'unaccented' targets in the prosodically intact sentences ($F(1,36) = 10.38$, $p < .005$). Moreover, there was a significant difference in the reverse direction, i.e. a response time advantage of originally unaccented over originally accented targets, in the prosodically manipulated sentences ($F(1,36) = 6.83$, $p < .02$). That is, when all three components of the prosodic contour signalled that accent would occur at the position where the target occurred, the target was responded to faster; and this was true whether the consistent prosody was applied to its original utterance or to its original utterance's pair.

This result allows us to dispose of the suggestion that prosodic transposition might interfere with listeners' prosodic processing. Instead, it is clear that what interfered most strongly with listeners' prosodic processing in the two preceding experiments was prosodic conflict. When one prosodic dimension was in conflict with the other two, listeners were unable to arrive at a consistent interpretation based on prosodic information. One effect of this was that significant accent effects disappeared.

However, the results from the prosodically manipulated conditions in Experiments 1 and 2, though they were both statistically insignificant, seem to differ. This might suggest that more sensitive experimentation could yet uncover differential contributions to the accent effect on the part of rhythm, pitch and intensity respectively. For the present, though, we may conclude with confidence that listeners' processing of prosody is not simply an additive evaluation of separate dimensions; the interaction between prosodic dimensions is of paramount importance. When the three dimensions rhythm, pitch and intensity agree, listeners exploit them efficiently and consistently. When they conflict, this exploitation is significantly impaired.

ACKNOWLEDGEMENTS

This research was supported by a grant from British Telecom to A. Cutler and D.G. Norris. Acknowledgement is made to the Director of Research of British Telecom for permission to publish this paper. The author is grateful to Jeffrey Bloom and Kim Silverman, without whose prosodic manipulation techniques these experiments would not have been possible, as well as to Steve Bartram, Sally Butterfield, John Williams and John Culling for technical assistance.

REFERENCES

- [1] Bloom, P.J. Use of dynamic programming for automatic synchronization of two similar speech signals. *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, 1984, 2.6.1-2.6.4.
- [2] Culler, A. Phoneme-monitoring reaction time as a function of preceding intonation contour. *Perception & Psychophysics*, 20, 1976, 55-60.
- [3] Culler, A. & Darwin, C.J. Phoneme-monitoring reaction time and preceding prosody: Effects of stop closure duration and of fundamental frequency. *Perception & Psychophysics*, 29, 1981, 217-224.
- [4] Norris, D.G. A computer-based programmable tachistoscope for non-programmers. *Behavior Research Methods, Instrumentation and Computers*, 16, 1984, 25-27.