# Timing Interference to Speech in Altered Listening Conditions

**Peter Howell** and **Stevie Sackin**
Department of Psychology, University College London

## Abstract

A theory is outlined that explains the disruption that occurs when auditory feedback is altered. The key part of the theory is that the number of, and relationship between, inputs to a timekeeper, operative during speech control, affects speech performance. The effects of alteration to auditory feedback depend on the extra input provided to the timekeeper. Different disruption is predicted for auditory feedback that is out of synchrony with other speech activity (e.g., delayed auditory feedback, DAF) compared with synchronous forms of altered feedback (e.g., frequency shifted feedback, FSF). Stimulus manipulations that can be made synchronously with speech are predicted to cause equivalent disruption to the synchronous form of altered feedback. Three experiments are reported. In all of them, subjects repeated a syllable at a fixed rate (Wing & Kristofferson, 1973). Overall timing variance was decomposed into the variance of a timekeeper (Cv) and the variance of a motor process (Mv). Experiment 1 validated Wing and Kristofferson's method for estimating Cv in a speech task by showing that only this variance component increased when subjects repeated syllables at different rates. Experiment 2 showed DAF increased Cv compared with when no altered sound occurred (Experiment 1) and compared with FSF. In Experiment 3, sections of the subject's output sequence were increased in amplitude. Subjects just heard this sound in one condition and made a duration decision about it in a second condition. When no response was made, results were like those with FSF. When a response was made, Cv increased at longer repetition periods. The findings that the principal effect of DAF, a duration decision and repetition period is on Cv whereas synchronous alterations that do not require a decision (amplitude increased sections where no response was made and FSF) do not affect Cv, support the hypothesis that the timekeeping process is affected by synchronized and asynchronized inputs in different ways.

## I. INTRODUCTION

Feedback control is one way of modeling on-line regulation of an action. Such models have been applied to speech control for which several potential sources of feedback are available (Fairbanks, 1955; Lee, 1950; Postma, 2000). Early models proposed that information about the identity of speech segments is used as feedback and that speakers take corrective action when the phonetic content of the actual speech output is discrepant from that intended (Fairbanks, 1955). Evidence that was originally regarded as strong support for this view was the disruption to speech performance that occurs when auditory feedback is altered experimentally (Lee, 1950). One such alteration is delayed auditory feedback (DAF), where speech is electronically delayed for a short time before it is replayed to the speaker. When DAF is administered, a speaker's accuracy and timing are affected dramatically (Lee, 1950). A monitoring account maintains that these effects arise because the experimentally-altered feedback, signals to the speaker that the speech was in error. The corrective action that is then taken results in disruption to speech control (Lee, 1950; Fairbanks, 1955).

Address correspondence to: Peter Howell, Department of Psychology, University College London, Gower St., London WC1E 6BT, ENGLAND. E-mail: p.howell@ucl.ac.uk.

Borden (1979) pointed out two fundamental problems in assuming auditory feedback is used to control speech: The checking and correcting operation cannot start until speech has been output. As processing the speech sounds and determining what action is necessary also take a significant amount of time (for instance, Marslen-Wilson and Tyler 1981 estimate recognition time for speech at about 200 ms), such feedback monitoring would lead to a slow speech output rate. Borden's second point was that speakers who are adventitiously deafened after they have learned a language, do not lose the ability to speak immediately, but that usually there is long-term degradation to speech control. The ability to speak immediately after hearing loss, suggests feedback is not necessary for on-line speech control (though the long-term effects may suggest feedback has a role in establishing and maintaining speech representations).

There are other problems, besides those raised by Borden, for this type of feedback view: The auditory system would need to supply the speaker with a veridical record of what was produced; otherwise, establishing if and what error has occurred with the intention of correcting it would not be possible. However, it is not clear that the representation of articulatory output provided by the auditory system is veridical of the intended message. This is because the auditory representation the speaker receives while speaking is affected by internal noise. The noise that is present then affects the information that can be recovered from the acoustic output. The main source of internal noise originates in vibrations of the articulatory structures that are transmitted to the cochlea through bone. This bone-conducted sound is delivered to the cochlea at about the same time as the acoustic output from the vocal tract. Bone-conducted sound during vocalization is loud enough to make its effects significant. von Bekesy (1960), for instance, estimated that bone- and air-conducted components are at approximately the same level. The air-borne sound contains sufficient information to decode a speaker's intention (other people listening to the speech understand the message). The bone-conducted sound, on the other hand, is dominated by the voice fundamental, formant structure is heavily attenuated and resonances of body structures extraneous to vocalization (such as the skull) affect this component (Howell & Powell, 1984). Consequently, the bone-conducted sound contains limited information about articulation. The degraded bone-conducted sound would also mask out the formant information in the air-conducted sound. Such masking would reduce the ability of a speaker to retrieve information about the articulation from the air-conducted feedback. This argument relies heavily on the evidence presented by Howell and Powell (1984). If future work confirms that the auditory feedback signal is restricted in the information it provides about articulation, models that assume feedback is used to compute a precise correction needed to obviate an error will need revision (Guenther, 2001; Neilson & Neilson, 1991, see the general discussion).

Another piece of evidence that the role of auditory feedback available over the short-term may have a more limited role than previously thought was given by Howell and Archer (1984). They reported an experiment in which a non-speech noise was substituted for the delayed signal under DAF conditions where the noise had the same intensity profile as the original speech. There was no difference in the time taken to read a list in this stimulus condition compared with that in a delayed speech condition. This suggests that any sound that stands in the same temporal relationship with the direct speech as does the DAF signal, will cause equivalent disruption to the delayed speech sound. Alteration to auditory feedback may not, then, have an effect on a monitor to correct segmental control, though it could still have a role in establishing and maintaining long-term speech representations (Borden, 1979).

The final piece of evidence about the role of feedback in short-term regulation of speech is based on changes that occur in voice intensity when altered auditory feedback (AAF) is

presented. When speaking in noise, voice level increases, whereas when concurrent speech level is increased, voice level decreases (Lane & Tranel, 1971). Howell (1990) reported that if DAF sounds are amplified, speakers increase vocal level. From this, it appears that delaying a person's speech creates a sound that is responded to as noise, rather than being processed as speech. A noise would be no use for indicating to a feedback mechanism whether adjustments to articulation are required (though, as noted earlier, auditory feedback could still have a role in establishing and maintaining long-term speech representations).

More recent models for the role of feedback have concentrated on (1) use of auditory feedback during development and recovery of speech control and (2) whether feedback can be used to check representations at suprasegmental (rather than segmental) levels. A representative example of each of these points of view is now described. The first approach can be illustrated in the work of Perkell and colleagues on normal speakers and speakers fitted with a cochlear implant. Perkell (1980) outlined a model in which feedback is used both for on-line control and monitoring overall speech output to make sure that it conforms with the norms of the language being spoken. The work reviewed earlier argues against auditory feedback having a role in on-line control, principally by questioning the interpretation of experiments that employ altered auditory feedback procedures. This would not necessarily undermine auditory feedback having a role in maintaining long-term representations. Also, a lot of this group's evidence is from patients fitted with cochlear implants and there are grounds, and evidence, for considering cochlear implant patients may be able to use auditory feedback for segmental control, irrespective of whether fluent speakers can use auditory feedback for this purpose.

Cochlear implant patients have no useful hearing before they receive an implant. The implanted electrode by-passes the peripheral auditory system and delivers sound direct to the auditory nerve. Auditory input can be presented to these patients once the electrode is switched on and any improved speech abilities investigated in a controlled way. Perkell, Guenther, Lane, Matthies, Vick and Zandipur (2001) reviewed their work on the role of auditory feedback in cochlear implant patients. They reported that restoration of hearing results in improved perception of vowels and of fricatives, production of vowels improved and this improvement occurred soon after patients received implants. Such findings support the view that these speakers may use feedback for control of speech. Note, though, that during vocal control, the information provided by the implanted electrode differs from that received by non-hearing impaired listeners. Bone-conducted sound is processed by the same peripheral mechanisms that transduce air-conducted sound (von Bekesy, 1960) and that have been damaged in these patients. Thus, when vocalizing, these patients would neither hear air-nor bone-conducted sound. The electrode only restores a representation of the air-conducted sound. The loss of bone-conducted sound during vocalization would prevent this sound from masking the airborne sound. The airborne component would then be a more useful source of information about what was articulated compared with normally-hearing individuals. The better information carried in the airborne source, relative to normally hearing speakers, would suggest that auditory feedback could be useful for control in these speakers (as Perkell's group has established).

An example of the view that auditory feedback is important for maintaining suprasegmental representations is reported by Natke and Kalveram (2001). They delivered frequency shifted feedback (FSF) shifted an octave down and measured any compensation speakers made. They found significant compensation on long, stressed vowels but not on short, unstressed ones. They argued that the compensatory response was based on a negative feedback mechanism. They also argued that the effects occur at a suprasegmental level since they were only observed on the long vowels. As voice pitch is available in both air and bone

signals, long-latency feedback mechanisms could extract it and it might then serve as a basis for suprasegmental control (i.e. in the way Natke and Kalveram, 2001, propose).

The aim of this article is to assess how AAF affects speech control in the short-term in normally-hearing individuals. The work does not address the issue whether auditory feedback is used outside these experimental procedures directly, nor how an internal model is established nor whether auditory feedback is used for monitoring prosodic aspects of speech. This has a different focus from the work with cochlear implantees where the auditory feedback is different from that received by normally hearing/listening speakers and from work investigating whether auditory feedback mechanisms have a role in suprasegmental control. The review of work on the involvement of a feedback mechanism in segmental control, offered above, suggests that in hearing individuals, the immediate effects of AAF do not arise because of interference with a feedback monitoring mechanism. The basic assumption behind the current tests is that altering sounds spectrally or temporally creates conditions that can lead to segregation (Bregman, 1990) of auditory feedback (principally bone-conducted when headphones are worn) from the altered sound. Depending on the alteration made, the altered sources that segregate will be asynchronous or synchronous with direct auditory feedback. The form of synchrony that arises when an alteration is made will determine the amount of disruption experienced.

Synchronous and asynchronous components that occur under different forms of AAF are considered first. Vocal output transmitted through bone will always remain in synchrony with articulation. With DAF, the delayed version of the speaker's voice is, then, asynchronous with this sound. In frequency shifted feedback (FSF), the voice is spectrally shifted with negligible delay (Howell, El-Yaniv & Powell, 1987). Bregman's (1990) work shows that such spectral manipulations (and, indeed, temporal alterations like those under DAF) lead to perceptibly distinct sound streams. So, in the case of FSF, the shifted sound source would lead to a separate, but synchronous, component to that which arises from any bone-conducted sound. The air-conducted component can be amplified independent of the bone-conducted sound and replayed, again with negligible delay. The difference in amplitude of these sound components creates conditions that would again favor segregation of the air-conducted from the bone-conducted sound (Bregman, 1990). The segregated sources would remain in synchrony (as with FSF).

The next stage in the argument concerns reasons for thinking that speaking at the same time as synchronous sound sources will be less disruptive on voice control than when speaking at the same time as asynchronous sound sources. Howell, Powell and Khan (1983) described several frequently-encountered situations that show asynchronous sounds are far more disruptive on speech control than synchronous sounds. This is demonstrated by considering two forms of song. Canons can be described as singing one song at the same time as another synchronous rhythm is heard. This is simple to perform, as shown by the fact that children are taught this form of song. One possible reason why tasks like canon singing are simple is that speaking or singing along with synchronous sounds reinforces the timing of the singer's own attempt, giving the listener a clearer sense of the beat that has to be followed (see Howell and Sackin, 2000, for supporting evidence). The second type of singing, hoquetus, is a mediaeval form in which different singers produce notes at the offset of the notes of a previous singer. So, the singer hears an asynchronous rhythm as well as that he or she produces. This form of song is difficult to perform relative to canon singing because of the presence of the asynchronous rhythm. According to the current hypothesis, extraneous synchronous and asynchronous rhythms create parallel situations to FSF or amplified feedback, and DAF, respectively. The effect of synchronous and asynchronous concurrent signals on voice control appears to be general (it is shown in AAF tasks and different forms of song). As this disruption is general, it suggests that the influences do not arise via the

operation of a feedback mechanism (it is unlikely that hoquetus singers treat other singers' notes as feedback about their own voice control that then gives rise to the observed disruption).

The next step is to consider what mechanism could explain differential disruption in situations where synchronous and asynchronous rhythms are heard, irrespective of whether these rhythms derive from a speaker's own speech or from an external event. As argued earlier, AAF manipulations can be regarded as transforming the speech task into one with additional rhythmic inputs (synchronous or asynchronous depending on the form of AAF). As the activity is serial, they could constitute input to a general-purpose timekeeping mechanism and the nature (synchronous or not) and number of inputs could then produce disruption through the timekeeper. Other serial activity would input to this same mechanism and give equivalent disruption depending, again, on the nature (relative to current input) and the number of such inputs.

Arguments for a central timekeeper that functions separately from motor processing activity have often been put forward in motor control since Lashley's (1951) seminal report on serial order behavior. (see for example Ivry, 1997; Vorberg & Wing, 1996; Wing & Kristofferson, 1973). Lashley argued that successive elements in a serial activity may be timed without reference to the peripheral motor events they give rise to. In particular, the completion of a motor element is not necessary for the generation of the next element in a sequence (as in a feedback process). One advantage of having a timekeeper that is independent of the execution of a particular act is that it can be used in a variety of tasks where timing control is needed.

Wing and Kristofferson (1973) established the properties of such a general-purpose timekeeper with data from a tapping task. Subjects started by tapping along with an isochronous metronome click. Once responses were entrained to the metronome's rate, the click was switched off and subjects continued the response sequence on their own. Variance associated with timekeeping processes (clock variance, Cv) was estimated separately from that associated with motor (motor variance, Mv) components. Wing and Kristofferson (1973) assumed that when a motor response deviated from its required position, two intervals were affected: If the response was ahead of its required position, the preceding interval would be short and the following interval would be long. Based on this, an estimate of 2Mv was obtained from the negative covariance between adjacent response intervals (lag one autocovariance). Cv and other residual components were then obtained by subtracting 2Mv from the total variance (Tv). The Cv estimate, unlike the Mv estimate, was not theoretically motivated. In further work, Wing (1980) validated that the residual provided an estimate of Cv by showing that as the length of the interval subjects were required to repeat (repetition interval) was increased, Cv also increased. This would be expected if keeping the time of long intervals is more difficult than keeping the time of short intervals and if difficulty is reflected in more variable responses. Mv, on the other hand, remained constant across repetition interval. Howell, Au-Yeung and Rustin (1997) have reported a similar validation to Wing (1980) for a task involving repeated movement of the lower lip (see also Hulstijn, Summers, van Lieshout, & Peters, 1992 for another application of the Wing-Kristofferson task to speech).

Work has also been conducted to establish the general properties of the timekeeper by examining whether Cv measures relate to other timekeeping operations. For instance, Ivry and Hazeltine (1995) examined production of specified intervals and, in separate tests with the same subjects, examined the relationship to subjects' perceptual time estimation ability. They reported a significant relationship. This relationship suggests that the mechanism has a general role in timing very different tasks.

The three experiments reported below test the hypothesis that AAF creates an ancillary sound that mainly affects the timekeeper (Cv) when speech is taking place concurrently. All the experiments use a speech version of the Wing-Kristofferson task and apply these authors' analysis procedure to estimate Cv and Mv. The speech version of the Wing-Kristofferson task involves repeating the syllable /bae/ at specified repetition intervals. The requirement to produce a single syllable with exact timing renders the task different to spontaneous speech. However, without these artificial task constraints, Cv cannot be estimated separately from Mv. Though there is this limitation, the same limitation applies whenever experimental techniques are used to study speech. Experiment 1 validates application of the Wing-Kristofferson method of obtaining Cv from residual variance after Mv has been extracted from Tv in the same way as in Wing's (1980) experiment discussed earlier. The second experiment tests whether synchronous and asynchronous forms of AAF reduce and increase Cv, respectively (Howell et al., 1983; Howell & Sackin, 2000). The delayed sound during DAF procedures is asynchronous relative to direct speech, so should lead to large increases in Cv (by analogy with the observations about performance disruption that occurs when activities are asynchronous, Howell et al., 1983). The second form of AAF tested is FSF in which the feedback is shifted in frequency with a negligible time delay. The psychoacoustic work described above shows two coincident spectrally-different streams of sound segregate (Bregman, 1990). Thus, in the case of FSF, the altered signal will separate from bone-conducted sound giving two synchronous inputs to the general-purpose timekeeper. In contrast to DAF, these should give a better beat so should not cause an increase in Cv. Experiment 3 investigates whether another way of changing synchronous input to the timekeeper (an intensity change) has similar effects to FSF. The effects of a secondary decision task on timekeeper operation were also examined in this experiment. Based on Ivry and Hazeltine (1995), a secondary task was chosen that involves a duration decision (i.e. time-based). It is assumed that the timekeeper is only sensitive to time-based decisions such as this (this is not tested explicitly). Consequently, Cv will only be affected when a duration decision is required.

## II. EXPERIMENT 1

The Wing-Kristofferson analysis is applied to data from a speaking task that requires the syllable /bae/ to be repeated. Total variance is decomposed into Cv and Mv components. The principal question is whether Cv, but not Mv, increases with repetition period (Wing, 1980) as validation that the Wing-Kristofferson procedure to obtain Cv applies to a speech task.

### A. Subjects

Eight adults (five males, three females) were employed. They had no history of speech or hearing disorder. They ranged in age from 26 to 34. The same subjects were used in Experiments 2 and 3 with half the subjects doing the experiments in one order, and the other half in reverse order (with the exception mentioned in the procedure for Experiment 3, below). Counter-balancing was done to avoid practice effects and to permit comparison across experiments (condition-specific practice for each type of trial was also given to avoid this problem). Subjects did each experiment on different days.

### B. Procedure

Subjects were told that the aim of the experiments was to establish the accuracy of articulatory timing when speaking a single CV syllable (/bae/) repeatedly at selected fixed rates. The syllable /bae/ was used because it is easy to say and its onset can be located reliably (the analyses are made from stop release). Subjects were instructed that on any trial the experimenter would play a recorded /bae/ at 70 dB SPL repeatedly at a particular rate.

Subjects were told that when they were ready, they should take a deep, but not excessive, breath (as though preparing to say a long sentence). They should then start their own production in time with the recording. The subjects were monitored to ensure they did not take a breath within a sequence. When they were going at the requisite rate the /bae/ used for entraining the speaker was switched off manually by the experimenter after the subject had responded to a minimum of five consecutive recorded /bae/s. The subject continued until either the experimenter stopped them or they stopped themselves because they had run out of breath. They were told not to take a breath in the middle of a sequence. The experimenter required the subject to continue the sequence for a minimum of 11 /bae/s. Four /bae/-/bae/ repetition periods were used (200, 400, 600 and 800 ms). The experiment started with practice at each repetition period until the experimenter judged that the subject was familiar with the task and could synchronize to the target at each rate. They then did the four different rates eight times each in a predetermined random order.

Subjects were tested individually in an Amplisilence sound-attenuated booth. The entrainment-/bae/s were played binaurally from a Toshiba laptop fitted with a Soundblaster 16 sound card. These sounds were relayed to Sennheiser HD480II headphones via a Fostex 6301B amplifier. Level of speech feedback after entrainment was set so that it was comfortable for listening (typically around 70dB SPL). Level was periodically checked. Speech was transduced with a Sennheiser K6 microphone and recorded on a DAT recorder.

## C. Analysis

The recordings were transferred to disk for analysis (48kHz sampling rate, 16-bit samples). The recordings were downsampled to 10kHz. /bae/-onsets were manually marked on 11 /bae/s in the phase after the entrainment sequence had been switched off starting at the onset of the first /bae/. Following Vorberg and Wing (1996), linear trends in the data were removed to ensure stationarity in the sequence. Inter-onset durations were calculated and Mv and Cv were computed from the algorithm given in Wing and Kristofferson (1973) on the ten intervals. The Wing-Kristofferson model only applies where the lag one autocorrelation, r, is bounded by $-0.5 < r < 0$ and some of the raw values lay outside these limits. So as not to bias the data by dropping these trials, intervals from each end of the series were dropped (a minimum of four intervals had to remain) from the original sequence, the truncated series was detrended and examined to see if it then fitted the Wing-Kristofferson model. The longest sequence that fitted was used in subsequent analyses. This allowed 98% of all trials to be included in the analysis. An analysis was also conducted to check this procedure does not affect the results. The trials where only the whole sequences fitted the model were used in these analyses. Analyses of data prepared in this manner produced equivalent results to those reported below.

## D. Results

Cv and Mv are plotted over repetition periods in Figure 1. A two-way repeated measures analysis of variance (ANOVA) was conducted with factors source of variance (Cv or Mv) and repetition period (200, 400, 600 or 800 ms repetition periods). The main effect of repetition period ($\underline{F}$ 3,21 = 13.3, $\underline{p}$ < 0.001) was significant. The Cv/Mv by repetition period interaction ($\underline{F}$ 3,21 = 5.60, $\underline{p}$ < 0.005) was also significant. Separate ANOVAs using either Cv or Mv alone showed Cv increased significantly as repetition period increased ($\underline{F}$ 3,21 = 13.3, $\underline{p}$ < 0.001) but Mv did not. The only repetition period that showed a significant increase in Cv over other intervals is 800 ms. The fact that significant differences between Cv and Mv occured at the longest interval is to be expected on the basis that Cv alone is affected when repetition period is lengthened.

### E. Discussion

The pattern of variance estimates (Cv and Mv) with change in duration of the repetition period is similar to that reported by Wing (1980) for a manual tapping task and by Howell et al. (1997) in a lip-tracking task. There is no change in Mv over repetition periods whereas Cv increases with the increase most apparent at the longest repetition period (800 ms). Wing (1980) argued that the selective increase in Cv with repetition period arises because of the greater difficulty controlling the timing of longer intervals. If Wing's (1980) reasoning is correct, the present results show that Cv provides an estimate of timekeeping processes in a task involving the speech articulators. Repetition periods of 600 and 800 ms are used in Experiments 2 and 3 to check whether feedback, intensity and decision-task manipulations affect this pattern. The repetition periods chosen (600 and 800ms) are in the region where Cv increases occur.

## III. EXPERIMENT 2

Speakers performed the Wing-Kristofferson task while listening to one of two forms of AAF (FSF and DAF). According to the hypothesis, both these types of AAF create an additional sound source as input to the timekeeper. The extra sound source in the case of FSF is effectively synchronous with speech and arises due to the spectral difference between the direct speech and its altered form. The additional sound source under DAF arises because of the temporal disparity between the direct and altered forms and so is asynchronous relative to the speech. According to the current hypothesis, asynchronous inputs to the timekeeper (as with DAF) cause more difficulty in performance than synchronous events. Consequently, the effect on Cv should be greater when DAF is presented than when FSF is presented. As DAF delay increases, the asynchrony between direct and delayed sources increases (in the experiment going from 66 through 133 to 200ms DAF-delay). More disruption should occur to the general purpose timekeeper as asynchrony increases. The effect on Cv and Mv while hearing each form of AAF is established, again, using the speech variant of the Wing-Kristofferson task.

### A. Participants

The same eight subjects were used as in the other experiments.

### B. Procedure

All conditions were performed as in Experiment 1 with the addition that one of the different forms of AAF was also heard. Besides this, the basic task was the same as in Experiment 1. The subjects were told to ignore the feedback and attempt to continue at the specified rate. On a trial involving a delayed sound, subjects heard standard DAF at one of three delays (66, 133 or 200 ms). (As argued in the introduction, in this, and all DAF experiments, speech is always transmitted through bone-conduction.) Subjects were tested at each DAF delay at repetition periods of 600 and 800 ms. Subjects received eight practice trials at each repetition period and DAF delay and then performed eight test trials at the same repetition period and DAF delay. The DAF delay and repetition period conditions were received in random order. The procedure for FSF was the same except that a time-synchronous, half-octave, downward frequency shift was fed back rather than a delayed sound.

The entrainment /bae/ sequence was played over a Toshiba laptop and Fostex monitor at the required repetition period, as in Experiment 1. Two Sennheiser K6 microphones were used to pick up the speech. One microphone supplied speech to a DAT recorder for use in the analyses. The other microphone output was relayed via a Quad microphone amplifier to the Digitech model studio 400 signal processor that produced the selected form of AAF. The Digitech output was played binaurally over Sennheiser HD480II headphones. The output

(set at 70 dB SPL) is at approximately normal conversational level so, according to von Bekesy's (1960) calculation the bone-conducted sound should also be roughly at this level too. The data were analyzed as in Experiment 1.

## C. Results

The results from the DAF conditions are shown in Figure 2. Mvs are plotted on the left and Cvs on the right. The axes are repetition period (abscissa) and variance (ordinate), as for Experiment 1, and the data points under the same DAF delay are connected together (identified by symbol). The results for the FSF condition are shown in Figure 3 in an equivalent way to the results of Experiment 1 (Cv and Mv estimates for each repetition period).

The Mvs in the DAF conditions were examined first with respect to whether DAF has a similar pattern to the results in normal listening conditions. Separate analyses were conducted for each DAF-delay and for each variance component to assess how DAF affects Mv and Cv relative to results on the same repetition periods in Experiment 1. For Mv, three two-way repeated measures ANOVAs were conducted. One factor was listening condition (the normal listening condition from Experiment 1 was always included and was compared with the selected DAF-delay condition, 66, 133 or 200 ms, from the current experiment). The second factor was repetition period (600 or 800ms). For the Mv measurement, no significant differences occurred between the normal listening condition and the 66 ms and 133 ms DAF-delays. The difference between normal listening and DAF was significant at 200 ms ($\underline{F}$ 1,7 = 11.4, $\underline{p}$ < 0.025) DAF-delay. These findings indicate that Mv increases under DAF over normal listening only at the longest delay. No interactions with listening condition were significant, so even with the most severe form of DAF, Mv appears to have the same pattern across repetition periods as in normal listening. Corresponding ANOVAs on Cvs showed normal listening differed from DAF at 133 ($\underline{F}$ 1,7 = 17.1, $\underline{p}$ < 0.005) and 200 ms ($\underline{F}$ 1,7 = 30.2, $\underline{p}$ < 0.001) delays. Differences across repetition periods were significant at all DAF-delays as main effects (66ms: $\underline{F}$ 1,7 = 11.50, $\underline{p}$ < 0.025; 133ms: $\underline{F}$ 1,7 = 8.5, $\underline{p}$ < 0.025; 200ms: $\underline{F}$ 1,7 = 29.00, $\underline{p}$ < 0.001). However, no interactions involving listening conditions were significant so there are only absolute differences between normal listening and DAF conditions

A three-way repeated measures analysis with factors DAF-delay condition (66, 133 or 200ms delay), variance source (Cv or Mv) and repetition period (600 or 800 ms) was next conducted to assess whether DAF-delay differentially affects Cvs and Mvs. DAF delay was significant ($\underline{F}$ 2,14 = 16.5, $\underline{p}$ < 0.001) showing DAF increases variances. These was also a difference between variance sources ($\underline{F}$ 1,7 = 60.0, $\underline{p}$ < 0.001) due to Cvs being greater than Mvs. Repetition period was significant ($\underline{F}$ 1,7 = 7.6, $\underline{p}$ < 0.05) with higher variances at the longer repetition period. The interaction of the latter factor with source of variance component shows higher variance at the longer repetition period. This is due to Cv increasing more over repetition periods than Mv does ($\underline{F}$ 1,7 = 4.10, $\underline{p}$ < 0.01). This result would be expected from Wing (1980) and Experiment 1. DAF-delay condition interacted with variance source ($\underline{F}$ 2,14 = 6.9, $\underline{p}$ < 0.01). This suggests that Cv and Mv increase at different rates with DAF-delay. Inspection of Figure 2 confirms this is most marked for Cv; Mvs exhibit less increase than Cvs (Mvs increase roughly three-fold over delays while Cvs increase more than five-fold).

A two-way ANOVA in which normal listening was compared with FSF (factor one) and repetition period (factor two) failed to reveal any significant differences. The equivalent two-way analysis on Cvs with normal listening and FSF as one factor and repetition period as a second factor showed a significant effect of repetition period ($\underline{F}$ 1,7 = 10.6, $\underline{p}$ < 0.025) but no further effects. The lack of an effect of FSF/normal listening as main effect or in

interaction shows that performance under FSF was not distinguishable from speech produced under normal listening conditions. As Cv increased in normal listening conditions over repetition period in Experiment 1, this might suggest that this occurs with FSF too. If so, it is surprising as Figure 3 appears to show little increase over repetition periods. This was investigated further in a two-way repeated measures ANOVA with factors variance source and repetition period on the data from the FSF condition alone. In this analysis, there was a main effects of Cv/Mv ($F$ 1,7 = 11.9, $p$ < 0.025) but the effect of repetition period was not significant. Interpretations based on effects that are not significant are problematic. Taking the two analyses together, the cautious conclusion would be that there is some attenuation of the increase in Cv over repetition periods (explaining why no Cv/Mv by repetition period interaction occurred when FSF alone was analyzed) but the attenuation is not detectable when Cv is compared across repetition periods between normal and FSF listening.

To summarize the findings, the pattern of Cv/Mv results over repetition periods shows that the global pattern of results under DAF is similar to what occurs under normal listening conditions (Mv is flat while Cv increases over repetition period). FSF, also shows no increase in Mv over repetition periods but, more surprisingly, little evidence for an increase in Cv over repetition periods. The other major finding is that Cv increases more than Mv as DAF-delay is increased.

## D. Discussion

When DAF was given, Mv showed less increase over repetition periods than Cv for these delays. Note that this general pattern, once again, validates the Wing-Kristofferson model for decomposing variance components (Wing, 1980). As well as the increase over repetition periods, Cvs increased more as DAF-delay was increased from 66, through 133 to 200 ms. The prediction that DAF should cause a marked increase in Cv with longer DAF-delays was confirmed. FSF produced a pattern of results in which Cv did not markedly increase as repetition period lengthened. The lack of increase in Cv under FSF at long repetition periods could be because this form of AAF is in synchrony with activity associated with speech. As argued in the introduction, the auditory feedback through bone and the FSF that is synchronous with speech, reinforce the timing of the direct speech giving the listener a clearer sense of the speech beat. This would help maintain the rate of the entrainment sequence leading to more precise control by the timekeeper (lower Cv).

The hypothesis that Cv does not increase at long repetition periods when the timekeeper has an input with a dominant beat is tested further in Experiment 3. In the introduction, it was argued that amplification can be considered as a form of time synchronous manipulation (Lane & Tranel, 1971). A better beat should arise when speech feedback is amplified so it should reduce Cv too (in this case, the louder the signal synchronized to speech activity the better sense of beat a speaker has available to maintain on-going timing).

## IV. EXPERIMENT 3

In Experiment 2, FSF (synchronous feedback) led to less of an increase in Cv at the long repetition period compared with Mv in the same listening condition and compared with the increase in Cv observed under DAF and under normal listening conditions. Amplitude alterations on the produced sequence of sounds were made in the current experiment as another form of synchronous sound alteration that should combine with sound sources that arise during vocalization, produce a better beat and prevent any increase in Cv at the longest repetition period. The experiment also included two conditions that varied what the subject had to do with the altered amplitude sections. A duration decision expected to tax the timekeeper (based on Ivry & Hazeltine's 1995 work) was selected as most likely to elicit

effects on Cv. In one condition, no response was made to the section where amplitude was altered (creating a condition similar to the FSF in Experiment 2). In the second condition, a duration response was made about the amplitude-altered section after the sequence was produced. These conditions were included to establish whether attention to the altered sound (such as when a response is required) leads to a Cv increase whereas, when attention to the sound is not required (as with FSF in Experiment 2), Cv does not increase. If a duration decision calls on the timekeeper's capacity, it would be revealed as increased Cv in conditions where timekeeping is most difficult (i.e. at long repetition periods).

### A. Participants

The same eight subjects were used as in the previous experiments.

### B. Procedure: Condition one (no response)

The basic experimental setup was the same for the two main conditions and was similar to that employed in Experiment 1. In condition one, some syllables were selected manually and amplified by 6dB by the experimenter as the speaker spoke. They heard the altered sound but did not make a judgement about it in this condition. The selection of syllables for amplification was made according to preset criteria (the remaining syllables were at the same level as in the previous experiments). Over all trials, the start of the first amplified section was either two or three /bae/s after the entrainment sequence stopped. This continued for two or three sounds. It was then presented at the standard level, again for two or three sounds. Finally, two or three more sounds were played at increased amplitude. This arrangement ensured that subjects did not know exactly when each of the two increased amplitude stretches in the sequences would start, or their duration. There were also equal numbers of trials that were the same (also with the same number of amplified sections containing two, or three amplified /bae/s) or different (same number of two/three and three/two amplified /bae/s). There were eight types of trial according to these criteria. These were presented eight times each in random order. Condition one was always received before condition two and condition three at the end to keep the subjects naïve as to the purpose of the amplified /bae/s. Subjects were given practice trials at the start. Only the two longest repetition periods were used (600ms and 800ms) because of the difficulty of the task.

### C. Procedure: Condition two (response)

The procedure for condition two was the same as for condition one, except that after they completed the interval production task, the subject was required to make a same or a different response about whether the two groups of amplified /bae/s had the same number of /bae/s or not. They were aware that they needed to make the decision before they started performing this condition. The experimenter gave feedback as to whether the subject was correct or incorrect after each trial (in actual fact, subjects were always correct).

### D. Procedure: Condition three

Condition three was a control included to check that subjects were able to perform the perceptual judgement in condition two accurately and how performance in condition two compared with a perception-only judgement. It did not involve production of sequences, only listening to, and making judgement about, recordings of the sequences obtained in condition two. This was done to ensure speakers did not maintain interval production performance by allowing duration decisions to be less accurate in that task. As duration decisions in condition two were perfect, this is superfluous. This condition is, however, described for completeness. The subjects listened to recordings of the sequences they had produced and did the same-different perceptual task alone. Performance on this task was (as in condition two) always correct. These results are not discussed further.

In conditions one and two, the entrainment /bae/ sequence was played over the Toshiba laptop and Fostex monitor at the required rate as in Experiment 1. Two Sennheiser K6 microphones were used to pick up the speaker's responses. One was led directly to the DAT recorder to be used later in the analysis. The other was routed via a Quad microphone amplifier to a Digitech model studio 400 signal processor. The output from the processor was split. Low amplitude white noise at about 60 dB SPL was added to the voice signal before the sound was fed back to the speaker to mask out the sounds of apparatus switching. This was played binaurally over a Sennheiser HD480II headset. The other output was recorded on a second channel of the DAT recorder (for use in the perceptual condition, condition three, and to check that the experimenter had made the alterations correctly).

### E. Results

The results are shown on the left of Figure 4 for the condition where the altered sound was heard but subjects did not judge the duration of the demarcated intervals and on the right for the condition where subjects did the duration judgement task after they completed the /bae/ synchronization task. Separate ANOVAs equivalent to those in Experiment 1 were conducted on each condition. In the condition where subjects made no response, there was a difference over repetition periods ($F_{1,7} = 17.2$, $p < 0.005$). Unlike Experiment 1, there was no interaction between Cv/Mv and repetition period. This suggests that the greater increase in Cv compared with Mv that occurred particularly at 800 ms repetition period in experiments 1 and 2 did not occur here.

A somewhat different pattern of results was found when subjects made a response to the altered sound. Cv/Mv ($F_{1,7} = 5.6$, $p = 0.05$), repetition period ($F_{1,7} = 7.3$, $p < 0.05$) and Cv/Mv by repetition period interaction ($F_{1,7} = 24.1$, $p < 0.005$) were all significant. The interaction shows Cv increased over repetition periods when a response was made (as expected from Experiment 1 and Wing, 1980). There was an interaction between response condition, source of variance component (Cv/Mv) and repetition period ($F_{1,7} = 8.2$, $p < 0.025$) when a three-way ANOVA was conducted with response condition as the extra factor. This shows that Cv increased at a different rate across repetition periods in the two response conditions (no increase in Cv over repetition periods when no response was made but an increase when a response was made).

### F. Discussion

Cv only increased over repetition periods if a response to sections increased in amplitude was required. This shows that this secondary decision affects the operation of the timekeeper. As only one task was used here, general disruption by any secondary task (rather than one specifically involving timing) cannot be definitely ruled out. However, Ivry and Hazeltine, (1995) found that performance on perceptual timing judgements correlates with variance in a tapping task. Given the very different nature of these tasks, it is difficult to see how there could be this relationship other than through a timing mechanism. The secondary task in the current experiment has similarities with that of Ivry and Hazeltine (1995) (insofar as a duration judgement is required). The main difference between Ivry and Hazeltine's work and the current is that in the latter the judgement is made concurrent with the interval production task rather than as a separate task. If it is accepted that the Ivry and Hazeltine's (1995) result showing correlations between the perceptual and production tasks operates through the timekeeping process, the interference from the secondary judgement task on Cv here would also seem to operate at the level of the timekeeper.

The next question considered is why there is no increase in Cv when subjects just listened to the sequences that had their amplitude altered. The results with FSF in Experiment 2, where again no increase in Cv over repetition occurred (though the lack of significance needs to be

treated cautiously), were explained by proposing that synchronized sound gives a better sense of beat to follow than occurs in normal listening conditions, leaving Cv less affected by repetition period. This explanation would apply in the condition where no response is given as the amplified sound gives a better beat and response requirements are the same as in the FSF condition. The general pattern of the results when a response has to be made to the altered sound is similar to that in Experiment 1 and as reported by Wing (1980) to validate the assumptions of the Wing-Kristofferson model. While the enhanced beat can remove the increase over repetition period, adding a duration decision adversely affects the timekeeper. The adverse effect of making a duration response is most evident in conditions where timekeeping is difficult (long repetition period). Note also that the effect of duration decisions on the timekeeper is consistent with the general role this mechanism is assumed to play (here, general to perception and production).

## V. GENERAL DISCUSSION

Experiment 1 validated the application of the Wing-Kristofferson task to speech and provided benchmark data against which to compare effects of AAF and additional response tasks. Experiment 2 showed that DAF has its principal effect on Cv whereas FSF did not lead to an increase in Cv at long repetition periods (note that this needs to be treated cautiously as it is based on finding no significant increase). The DAF result was predicted from the hypothesis that the DAF signal is asynchronous with direct speech and asynchronous inputs markedly affect timekeeper operations. The FSF effect was predicted on the basis that the altered sound is in synchrony with the ongoing speech sound input to the timekeeper and two synchronized inputs give the timekeeper a better sense of beat that aids (or prevents degraded) performance. Experiment 3 tested the effect of enhancing the beat on timekeeper operation further by altering the amplitude of sections of speech output. This experiment also included conditions where subjects were, and were not, required to make a response about the sections with higher amplitude (enhanced beat). In the condition in which an amplitude alteration was heard but not responded to, Cv showed no increase at the long repetition period as occurred with FSF. However, Cv did increase when the same sounds were heard when a duration judgement was required about the amplified section. This suggests the Cv increase emerges when a duration decision is required because the difficulty faced by the timekeeper is enhanced.

The implications of the results are considered for the role of AAF experiments for speech control. The first topic considered is whether a case can be made that temporal alterations lead to temporal disruption of speech and non-temporal alterations lead to non-temporal changes in speech output. The evidence on FSF appears to speak against the second proposition. In conditions approximating closer to normal speaking conditions, at first sight, there is evidence that the spectral alteration of FSF does not produce any noticeable change in timing (for instance overall average speech rate appears normal). However, Howell and Sackin (2000) looked at FSF on sentence material and found that timing variability around specified segmentation points is affected by this manipulation. In particular, FSF reduced timing variability possibly by enhancing the direct beat. Thus, as spectral alterations during FSF lead to significant effects on timing control, the notion that only timing alterations cause speech-timing changes on vocal output cannot be sustained.

This leads on to the second issue, whether AAF experiments provide support for feedback monitoring models. The current work does not necessarily rule out auditory feedback having a role in maintaining internal long-term models for the speaker's language (Perkell, 1980). Auditory feedback might also have a role in control of segmental aspects in speech produced outside AAF procedures, depending on the position taken about whether the auditory feedback is reflexive of the speaker's intention or not. Currently there is only one study that

suggests auditory feedback is not reflexive of speech output (Howell & Powell, 1984). This suggests that the cautious approach would be to not definitely rule out auditory feedback having a role in maintaining gross segmental information. Having said this, a model that do not require reflexivity (Howell, in press) merits brief discussion. Examples of each of these two types of model are considered starting with two models that require reflexivity, Neilson and Neilson's (1991) adaptive model theory (AMT) and Guenther's (2001) DIVA model.

Neilson and Neilson's (1991) AMT theory has a controlled dynamic system driven by an adaptive controller. The adaptive controller transforms motor commands into sensory events. The adaptive controller has to have access to the speech-output "solution" that is obtained by an inverse dynamics process applied on sensory feedback. It would not be possible to determine whether and what correction is necessary if the sensory signal is non-reflexive of speech output. In Guenther's (2001) model, auditory targets are projected from premotor cortical areas to the posterior superior temporal gyrus where they are compared to incoming auditory information via primary auditory cortex. Any difference represents an error signal that is mapped through the cerebellum and the auditory error signal indicates a change is required to the motor velocity signal that controls the articulators to zero the error. Again if the auditory feedback is not sufficiently reflexive of speech output, information about segment articulation would not be veridical and could potentially even lead to incorrect corrections.

Howell (in press) has offered a model where the altered sound inputs to the timekeeper and how this causes disruption rather than the feedback from the altered sound continues to be used by a monitor for feedback control. The principal advantage of this model is that it avoids the non-reflexivity problem. This interpretation suggests alteration to auditory feedback creates an artificial speaking situation. This does not necessarily rule out a role for auditory feedback in segmental control in normal speaking situations. Howell's (in press) model circumvents the reflexivity problem by proposing cerebellar mechanisms give an error signal that only arises when timing problems occur. This alert acts as an all or none signal given the sole role of slowing speech rather than segmental correction. Loss of hearing just leads to one less input to the timekeeper and the timekeeper is not adversely affected by removal of this source of input. Initiation of a subsequent sound once one sound has finished does not depend on the results of processing sound back through an auditory feedback loop. As this rate-limiting step in speech control is removed, there is no problem in accounting for the rapidity at which speech can be produced. It would not matter, then, whether auditory feedback of the voice presents a veridical representation of what was said; it will only depend on the timing of the altered sound in relation to other timekeeper inputs. Thus, an altered sound that has the same timing as DAF speech would offer the same serial input to the timekeeper and produce equivalent disruption (Howell & Archer, 1984). Finally, manipulations that transform speech into noise that has lost its association with the original speech by being delayed (Howell, 1990) would be effective because of the asynchronous input they provide, not because the sounds were originally derived from speech. The EXPLAN model has its limitations. For instance, it does not address the issue about how long-term representations are established. Establishing how degraded auditory feedback is of speech output is a topic that merits further attention as it features in many monitoring models.

## Acknowledgments

# References

von Bekesy, G. Experiments in hearing. New York: McGraw-Hill; 1960.

Borden GJ. An interpretation of research on feedback interruption in speech. Brain Lang. 1979; 7:307–319. [PubMed: 455050]

Bregman, AS. Auditory Scene Analysis: The Perceptual Organization of Sound. Cambridge MA: MIT Press; 1990.

Fairbanks G. Selective vocal effects of delayed auditory feedback. J. Speech Hear. Dis. 1955; 20:335–348. [PubMed: 13272227]

Fowler CA. Coarticulation and theories of extrinsic timing. J. Phon. 1980; 8:113–133.

Guenther, FH. Neural modeling of speech production. In: Maassen, B.; Hulstijn, W.; Kent, R.; Peters, HFM.; van Lieshout, PHMM., editors. Speech motor control in normal and disordered speech. Nijmegen: Uttgeverij Vantilt; 2001. p. 12-15.Nijmegen

Helmuth LL, Ivry RB. When two hands are better than one: Reduced timing variability during bimanual movements. J. Exp. Psy: Human Percept. Perf. 1998; 22:278–293. [PubMed: 8934844]

Howell P. Changes in Voice Level Caused by Several Forms of Altered Feedback in Normal Speakers and Stutterers. Lang. Speech. 1990; 33:325–338. [PubMed: 2133911]

Howell, P. The EXPLAN theory of fluency control applied to the Treatment of Stuttering by Altered Feedback and Operant Procedures. In: Fava, E., editor. Current Issues in Linguistic Theory. Pathology and therapy of speech disorders. Amsterdam: John Benjamins; in press

Howell P, Archer A. Susceptibility to the Effects of Delayed Auditory Feedback. Percept. Psychophy. 1984; 36:296–302. [PubMed: 6522222]

Howell, P.; Au-Yeung, J.; Rustin, L. Clock and motor variances in lip-tracking: A comparison between children who stutter and those who do not. In: Hulstijn, W.; Peters, HFM.; van Lieshout, PHHM., editors. Speech Production: Motor Control, Brain Research and Fluency Disorders. Amsterdam: Elsevier; 1997. p. 573-578.

Howell, P.; El-Yaniv, N.; Powell, DJ. Factors affecting fluency in stutterers when speaking under altered auditory feedback. In: Peters, H.; Hulstijn, W., editors. Speech Motor Dynamics in Stuttering. New York: Springer Press; 1987. p. 361-369.

Howell P, Powell DJ. Hearing your Voice through Bone and Air: Implications for Explanations of Stuttering Behaviour from Studies of Normal Speakers. J. Fluency Dis. 1984; 9:247–264.

Howell P, Powell DJ, Khan I. Amplitude contour of the delayed signal and interference in delayed auditory feedback tasks. J. Exp. Psy: Human Percept. Perf. 1983; 9:772–784.

Howell P, Sackin S. Speech rate manipulation and its effects on fluency reversal in children who stutter. J. Devel. Phys. Disab. 2000; 12:291–315.

Hulstijn W, Summers JJ, van Lieshout PHM, Peters HFM. Timing in finger tapping and speech: A comparison between stutterers and fluent speakers. Hum. Movement Sci. 1992; 11:113–124.

Ivry, R. Cerebellar timing systems. In: Schmahmann, J., editor. The Cerebellum and Cognition. San Diego: Academic Press; 1997. p. 555-573.

Ivry RB, Hazeltine RE. Perception and production of temporal intervals across a range of durations: Evidence for a common timing mechanism. J.. Exp. Psy: Human Percept. Perf. 1995; 21:3–18. [PubMed: 7707031]

Lane H, Tranel B. The Lombard sign and the role of hearing in speech. J. Speech Hearing Res. 1971; 14:677–709.

Lashley, KS. The problem of serial order in behavior. In: Jeffress, LA., editor. Cerebral mechanisms in behavior. New York: John Wiley; 1951.

Lee BS. Effects of delayed speech feedback. J. Acoust. Soc. Am. 1950; 22:824–826.

Marslen-Wilson WD, Tyler LK. Central processes in speech understanding. Phil. Trans. Roy. Soc. Lond. Series B. 1981; 259:297–313.

Natke U, Kalveram KT. Effects of frequency-shifted auditory feedback on fundamental frequency of long and short stressed and unstressed syllables. J. Speech, Lang. Hear. Res. 2001; 44:577–584. [PubMed: 11407562]

Neilson, MD.; Neilson, PD. Adaptive model theory of speech motor control and stuttering. In: Peters, HGM.; Hulstijn, W.; Starkweather, CW., editors. Speech Motor Control and Stuttering. Amsterdam: Elsevier; 1991. p. 149-156.

Perkell, J. Phonetic features and the physiology of speech production. In: Butterworth, B., editor. Language Production. Vol. 1. London: Academic Press; 1980. p. 337-372.

Perkell, J.; Guenther, F.; Lane, H.; Matthies, M.; Vick, J.; Zandipur, M. Planning and auditory feedback in speech production. In: Maassen, B.; Hulstijn, W.; Kent, R.; Peters, HFM.; van Lieshout, PHMM., editors. Speech motor control in normal and disordered speech. Nijmegen: Uttgeverij Vantilt; 2001. p. 5-11.

Postma A. Detection of errors during speech production: A review of speech monitoring models. Cognition. 2000; 77:97–131. [PubMed: 10986364]

Vorberg, D.; Wing, A. Handbook of Perception and Action. Vol. 2. London: Academic Press; 1996. Modeling variability and dependence in timing; p. 181-262.

Wing, AM. The long and the short of timing in response sequences. In: Stelmach, GE.; Requin, J., editors. Tutorials in Motor Behavior. Amsterdam: North Holland; 1980. p. 469-486.

Wing AM, Kristofferson AB. Response delays and the timing of discrete motor responses. Percept. Psychophys. 1973; 14:5–12.
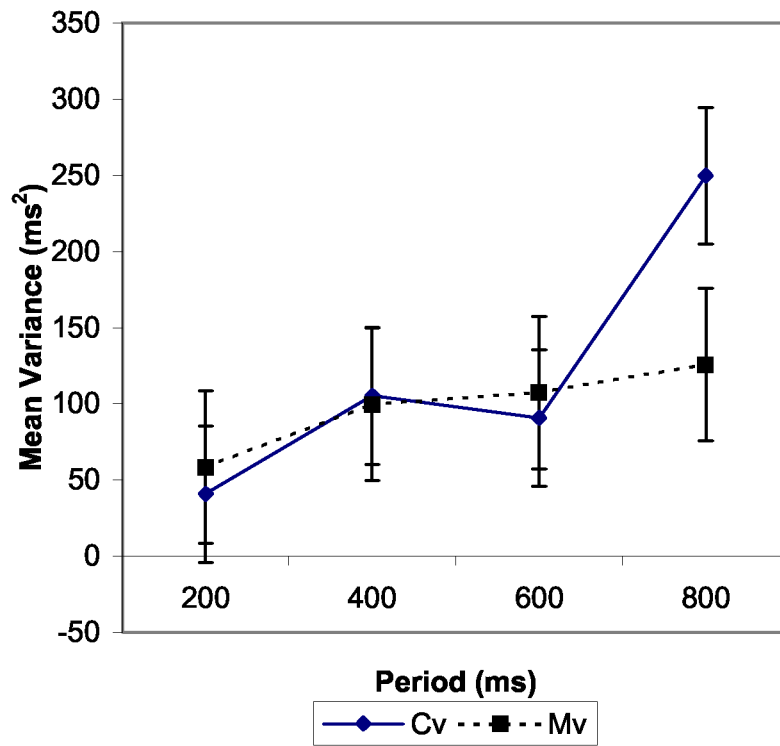
**Figure 1.**
Motor (dashed line) and clock (solid line) variances (ordinate) for repeating the syllable /bae/ at different periods (abscissa). Periods used were 200, 400, 600 and 800 ms.
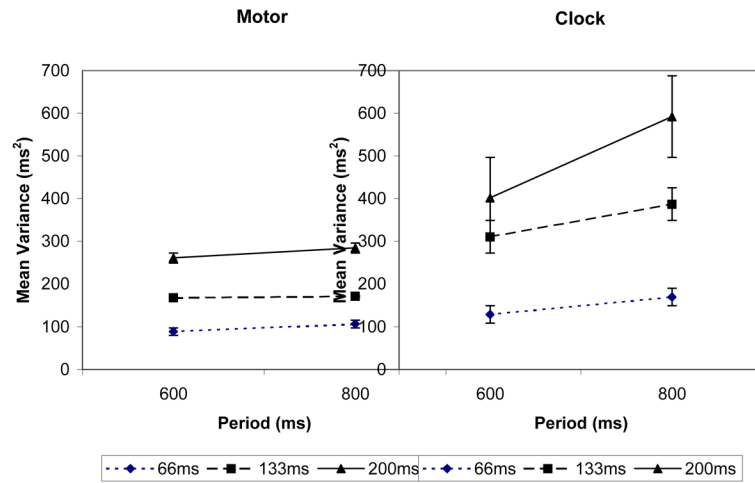
**Figure 2.**
Motor and clock variances against period (600 and 800 ms) for the delayed auditory feedback conditions of experiment 2. Motor variances are shown on the left and clock variances on the right. DAF delay of the points connected together was 66, 133 or 200 and the delay used for connected points can be identified from the symbol in the caption.
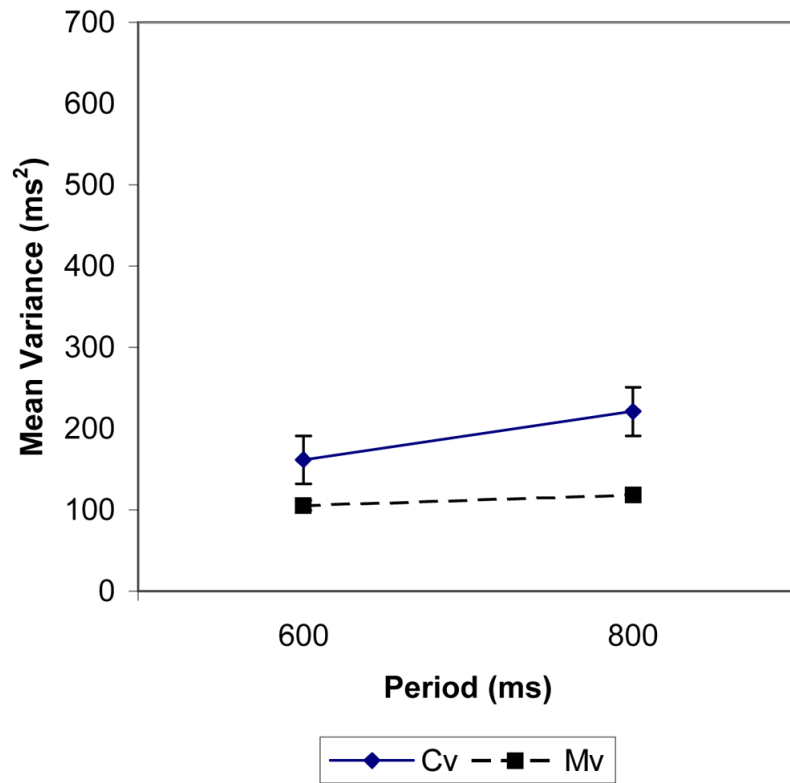
**Figure 3.**
Motor and clock variances against period (600 and 800 ms) for the frequency shifted condition of experiment 2.
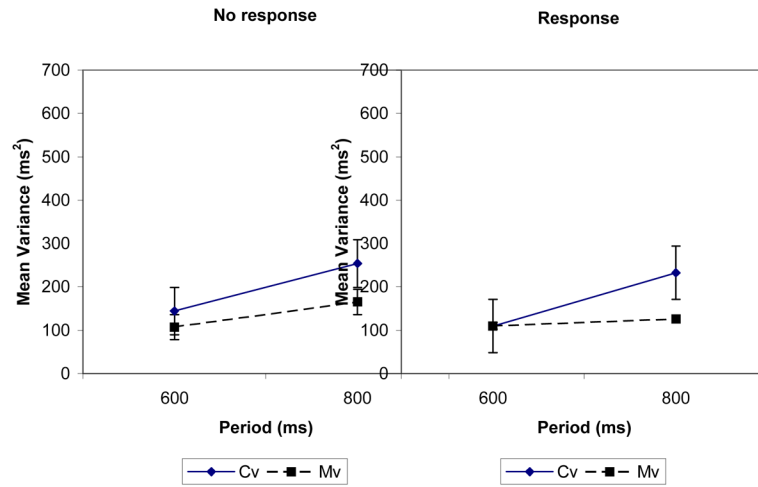
**Figure 4.**
Motor and clock variances against period for experiment 3. Periods used were 600 and 800 ms. The results for the no response condition are on the left and the response condition on the right.