

A Computer Implementation of Psychoacoustic Grouping Rules

Daniel P. W. Ellis

Perceptual Computing, MIT Media Lab E15-401B, Cambridge MA 02139
dpwe@media.mit.edu

Abstract

We are building a computer model of sound organization and understanding in human listeners. In particular, we would like to be able to detect and locate acoustic events that will be perceived as separate objects. Our model aims to duplicate this aspect of the auditory system, although the level of correspondence is speculative given our current state of knowledge.

We describe an implementation of grouping rules corresponding to the psychoacoustic cues of harmonicity, common onset, continuity and proximity [1]. We increase the system's robustness by adding a second layer of grouping that looks for corroboration between primary groupings. We believe that such a system of repeated hierarchic grouping is critical for the successful modeling of auditory functions.

1. Introduction

This paper describes a computer model of part of the human auditory system, specifically the process by which disparate acoustic energy incident upon the ears is 'organized' into a small number of separately-identified real-world sound sources. This problem of *source separation* is analogous to the segmentation problem in vision: the raw information available to the perceptual system is the combination of several essentially independent sources; the most effective way to process this information (avoiding combinatorial explosion) is to isolate the contribution of each source and deal with it separately; therefore, the first stage of processing must be concerned with identifying and grouping sets of information by source.

When energy in disjoint frequency bands is perceived as arising from a single source, the separate bands are said to be *fused*. Consider for example a complex periodic tone: The listener usually hears a single tone with pitch related to the common period of the harmonics, and 'quality' depending on the harmonics' amplitudes. Psychoacoustic experiments have resulted in a set of empirical rules to predict how combinations of simple sine stimuli will be organized by listeners [1].

Signal processing of sound has typically been limited in its ability to cope with interfering mixtures because of the lack of a 'segmenting' front-end. This is particularly evident in speech recognizers that can only function when given clean, isolated input. Lately, several researchers have been investigating approaches based on these rules of auditory organization: Cooke [2], Brown [3] and Okuno *et al* [4] have focused on the problem of enhancing speech amid interference, whereas Mellinger [5] and Kashino [6] have considered the related problem of separating the different melodic lines in polyphonic music.

The current project is motivated by a belief that one critical aspect of any successful model of human auditory grouping is the simultaneous use of a *range of different* cues. Moreover, the process of combining these cues is itself a central and unique facet of such systems, perhaps more important than the details of the individual cues themselves, making the auditory perception system a 'society' in the sense of Minsky [7]. Therefore we set out to build a selection of simple grouping schemes in order to be able to experiment with the problem of combining their outputs.

While the system we describe has an essentially ad-hoc and problem-specific structure, it is worth noting its similarity, both in outline and to some extent in detail, to the sound-understanding blackboard systems [8], an axis of abstraction we hope to investigate and develop in later work.

Although it would certainly be useful to endow machines with the abilities to organize and interpret sound that are possessed by people, our primary motivation has been to understand the human prototype rather than solve particular problems in automatic processing of signals. We chose to pursue this by building a functional model; this approach has several advantages as a method of testing the validity and interaction of complex theories. It may also end up performing some useful processing.

The next section gives an overview of the current project, and sections 3, 4 and 5 describe in more detail the representation, primary grouping and secondary grouping respectively. Section 6 introduces some preliminary results of applying the system to real sounds, and section 7 concludes with a brief discussion of our planned developments.

2. Overview of the system

Figure 1 shows how the current work fits into the context of a complete machine audition system. At the bottom of the diagram, real-world events generate sound, which is analyzed into a time-frequency representation by a filterbank approximating the function of the cochlea. The output of the filterbank is represented as time-frequency contours (or *tracks*) lying along energy maxima, as described in section 3. This representation is then subject to grouping into objects: The two rounded boxes labeled “primary grouping” and “secondary grouping” correspond to the work described in this paper.

The first box, primary grouping, is made up of rules for identifying the basic psychoacoustic cues such as harmonicity, common onset and proximity, and is described in more detail in section 4. These rules each produce multiple groups of time-frequency tracks.

These groupings are called *part-objects* in the diagram to underline their preliminary nature. They are fed to the secondary grouping stage which performs somewhat differ-

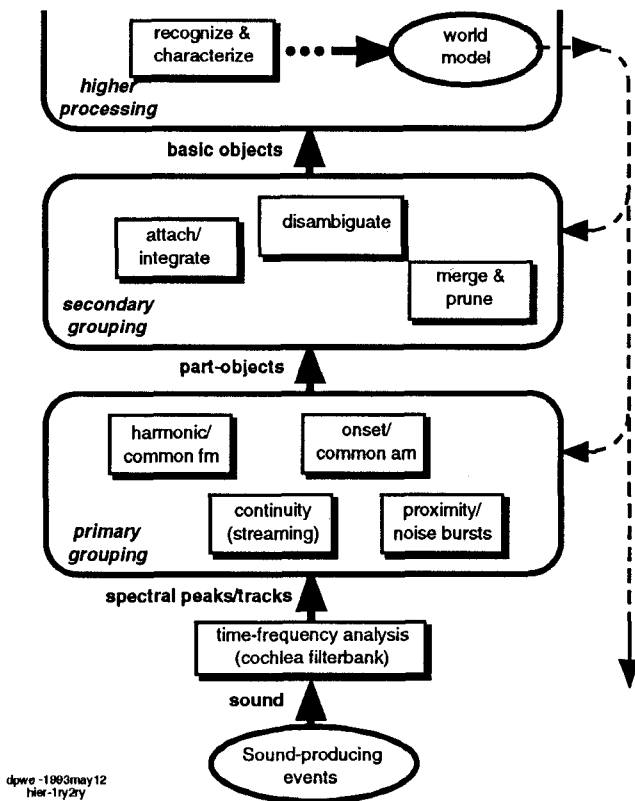
ent operations to produce more robust track groupings by integrating the results of the primary rules. These operations include looking for duplicate and (near) subsets among the groups and reducing these down to a smaller set of canonical groupings.

The pruned and non-redundant track groups produced by this second stage are the *basic objects* of the diagram, since they are very likely to be perceived as individual entities, yet they are still relatively simple in structure. In particular, they have been created only by consideration of general principles of the structure of real-world sounds, but without any knowledge or memories of sounds they might resemble or evoke. This kind of acquired or knowledge-based processing would be accomplished by the next layer, labeled “higher processing” in the diagram, and deliberately drawn as incomplete since we do not have strong ideas about how far upwards this block diagram should really extend. The dotted arrows pointing down from the “world model” indicate the possibility of ‘goal-driven’ processing that would guide searches to reveal such details based on expectations.

3. A representation for sound

The underlying representation used for the processing, the Constant-Q Sinewave Model [9] based upon the Sinusoid Transform [10]. We review the process very swiftly here. The continuous variation in air pressure as a function of time is converted into a set of discrete contours or ‘tracks’ via three stages: first, the acoustic energy is calculated as a function of time and frequency by passing the signal through a constant-Q filterbank. At each time-instant the magnitude and frequency of the local maxima of the spectrum are recorded; these are matched between successive time frames to create tracks. These tracks comprise a set of discrete pairs of frequency and magnitude functions that represent the original sound. Subsets of the tracks may be simply resynthesized via sinewave oscillators achieving a high perceptual similarity to the original sound. The key properties of the representation for the current application are that, excluding spectral collisions, each track represents energy of only a single source, and yet has meaningful attributes such as frequency modulation rate and magnitude variance.

Figure 2 shows the graphical format we use to display this bottom-level analysis, in this case the sound of a solo clarinet corrupted by the sound of a can dropping onto a hard surface. Time goes from left to right; this example lasts a little under a second as indicated by the scale at the bottom. The top panel shows the envelope of the magnitude waveform. The lower panel has logarithmic frequency as its vertical axis and shows the energy output of the filterbank as shades of gray. We see the first five or six clarinet harmonics separately resolved as in a narrowband spectrogram, and the broadband vertical structures of the can impacts around $t = 1.04, 1.16, 1.33$ and 1.47 s.



dpwe - 1989may12
hier-1ty2ty

Figure 1: Block diagram of a ‘complete’ sound processing system, including the primary and secondary grouping stages described in this paper.

The lines drawn over the filterbank output, particularly visible along the centers of the resolved harmonics, are the frequency contours of the tracks formed by the analysis described above. The can noise is represented as a number of short-duration, haphazardly scattered sinusoids with uncorrelated modulation.

4. Primary grouping rules

The first stage of grouping consists of more-or-less direct implementations of the cues known to us from psychoacoustics. Each rule is written to generate several answers when grouping is ambiguous. Thus the overall strategy is to generate a large number of track groups which will, with high confidence, contain the desired groups among many others, and then prune poor-quality and duplicate groups at a later stage. We consider each rule in turn.

4.1 Harmonicity rule

The grouping of acoustic energy by harmonicity derives from the real-world observation that sound is often generated by near-periodic mechanisms — the oscillation of the human vocal folds being a particularly important example. A Fourier analysis will represent a periodic signal by a series of distinct harmonics, at least for the region of the spectrum over which the analysis bandwidth is less than the spacing between the harmonics. These harmonics, which occur at integer multiples of the repetition frequency or *fundamental*, correspond to the Fourier series expansion of the periodic signal. The harmonicity rule finds regions of energy that show this pattern.

It is clear that the track representation is particularly suitable for this purpose, since each resolved harmonic will typically result in a single, isolated track. The harmonicity rule considers all tracks that overlap in time with the seed track, and that are close enough in frequency to possibly be one of the harmonics resolvable by our filter bank (the constant Q nature limits this; harmonics above the sixth or seventh will not be resolved).

For each candidate member of the harmonic group, the frequency ratio to the seed track is calculated for each sample point during their time overlap. In order to be added to the group, the ratio has to be close to integer, and close to constant over the track duration.

4.2 Common onset rule

The ‘ecological’ basis for the common onset rule (i.e. the regularity of the real world that makes it useful) is that if a particular physical process generates energy in various frequency bands, it is likely that energy will start in each of those bands at the same moment. This would seem to be simple to

detect among the tracks, however it is necessary to include some intelligent tolerance of asynchrony to accommodate both phase distortion of the channel, and the intrinsic time uncertainty of the narrow, low-frequency filter channels. Onset groups are iteratively extended from a low frequency seed, so a single seed may give rise to a number of plausible onset groups.

4.3 Continuity rule

This rule accomplishes grouping of tracks across short time gaps. This may be necessary due to energy modulation in the source, or can be viewed as a pragmatic measure to ‘clean up’ the output of the track-formation stage, by compensating for occasions when that stage ‘lost track’ of particular energy regions. We can organize the latter case into a continuity group (or *metatrack*) that in many ways behaves like a single, deeply modulated track.

The actual implementation is very straightforward; a time-frequency region just ahead of the end of the seed is searched for onsets of possible continuations. The continuation is terminated when energy of the tracks being added falls to some threshold below the average magnitude of the group.

4.4 Proximity (noise) rule

While harmonic complexes and resonant (formant) bursts are represented very successfully as sums of modulated sinusoids, such a transformation is less obviously appropriate for wideband, sustained ‘noise’ such as speech sibilants. These analyze into a large number of densely packed, uncorrelated tracks. We would like some method to group such tracks into the single perceptual element to which they correspond. A smoothed spectrum is obtained by convolving the spectral magnitude function at a particular instant with a Gaussian kernel; a candidate noise band is then defined as a range in this smoothed spectrum that lies within a magnitude threshold of its average value for at least a certain bandwidth. This search is then repeated at successive time instants; if highly-overlapped frequency bands are found at several adjacent times, a complete time and frequency range for the noise region is defined. Any tracks lying in this region of appropriate energy and duration are recruited to the corresponding noise group.

4.5 Amplitude-modulation groups

The last major monaural psychoacoustic grouping cue is common amplitude modulation, the association of energy in different frequency bands that exhibits synchronized fluctuations. Such a cue would help in collecting the separate formant trains of a single voice (since they will all be

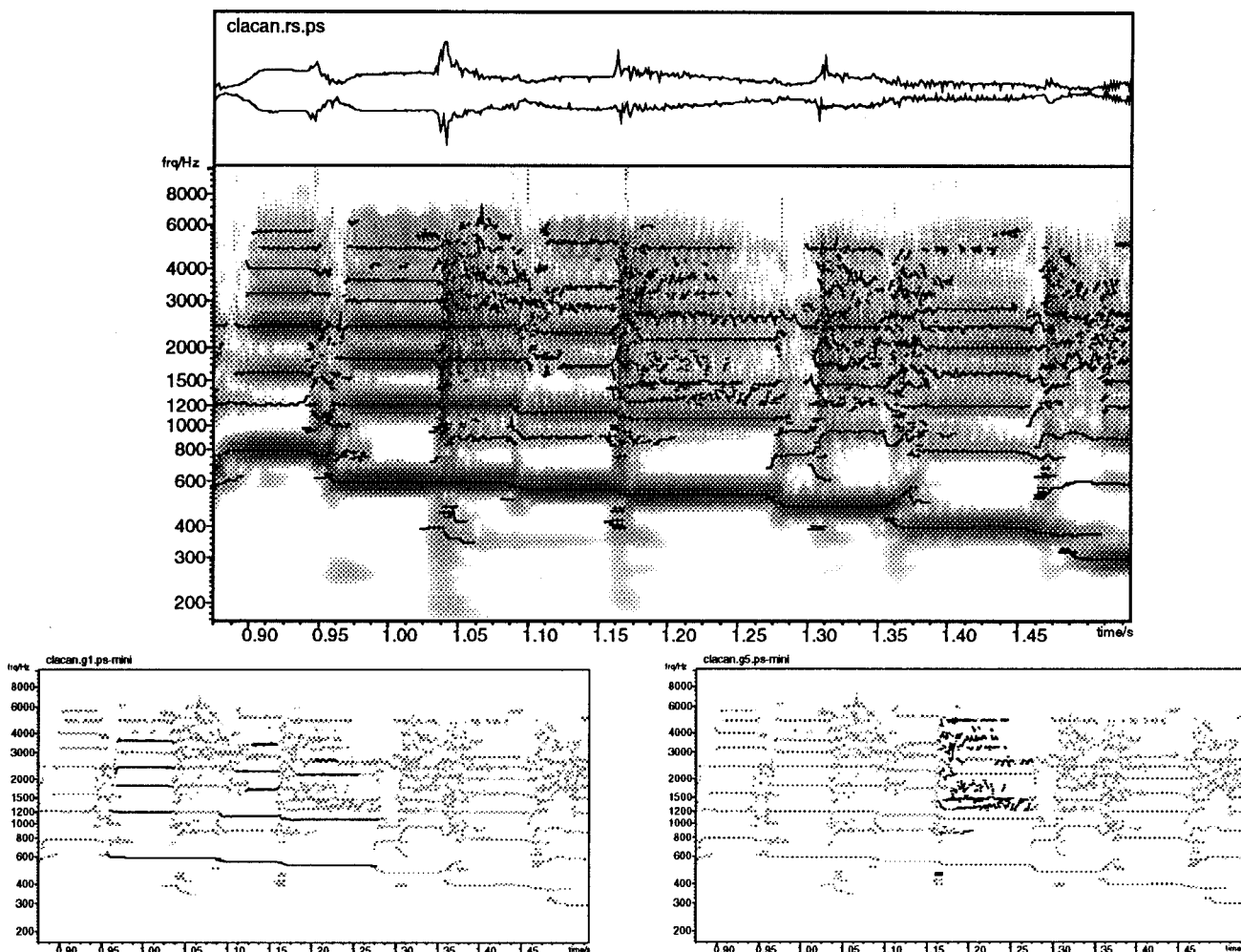


Figure 2: The top image shows the analysis of the original mixture, a clarinet melody with repeated intrusions of a can hitting a hard surface. The lower images show two different groups formed by the system as the emphasized tracks.

amplitude-modulated by the pitch pulses), and, with some modification, could help attach the formants to the resolved harmonics. We have yet to implement such a rule. This should be natural and straightforward - in a manner similar to the harmonic grouping, pairs of tracks (or metatracks) could be scored for amplitude modulation correlation during their common time support.

5. Secondary grouping

We have now described in detail our implementation of the primary layer of group-formation, whereby psychoacoustic cues are used to generate a large number of groups of tracks showing various degrees of association according to the cues. The groups themselves are not yet useful for interpreting or explaining the sound; there are too many of them, there are many redundant groups, and some of them have only a low confidence score. To produce a small number of less redundant, high confidence groups, we apply

a second stage of grouping which aims to prune and integrate the results of the primary rules.

5.1 Pruning

All the groups generated by a particular primary rule are sorted into one list, according to total energy. Then each group is compared to all the groups below it in the list; if a lower group is a proper subset of the larger group, or if the energy in their difference is less than some threshold, the smaller group is deemed a (near) subset which is adequately represented by the larger, and it is removed from further consideration. This typically effects a 10:1 or greater reduction in the number of track groups.

5.2 Correlation

Pruning is only applied within a particular rule since the existence of highly-overlapped groups *between* different

rules is very important evidence of genuine coherence for that group. The next stage of processing searches for this coherence by looking for highly-overlapped pairs. If such a pair has energy in their common tracks above some threshold of their total energy, a new 'supergroup' is formed by merging their component tracks.

5.3 Closure

In order to form robust groups that are not harmonically based we need some way of combining the results of continuity and proximity (noise) rules. Currently we merge together all such rules that intersect with principal onset groups. Then all non-harmonic groups are ranked by total energy and pruned as described above; those with the highest energy are carried forward as higher-level objects.

6. Results

The lower panels in figure 2 show two typical examples of the groups found by the scheme. The first identifies a harmonic cluster comprising clarinet notes, which results from the correspondence between onset and harmonicity groups. The second shows an individual can impact, which is the closure between an onset group and continuity and noise groups.

7. Conclusions and future work

We emphasize that this is preliminary work, and there are many outstanding issues we are anxious to address.

At the lowest level, there are still questions about the adequacy of the underlying track representation. The precise nature of the filterbank response is highly idealized (compared to physiological data) and is essentially linear; addressing either of these could have far-reaching implications for subsequent processing.

We would also like to provide for top-down processing, where the data provided by the early stages can be changed in response to higher level inferences and deductions. The subtraction of components in the case of an implicit spectral collision is a specific although difficult example in this area.

Considering the primary grouping rules, we have noted some shortcomings in the description above. Most pressing are the need for a common-AM grouping rule, and further refinements or innovations in a 'noise region' grouping rule.

For the secondary grouping stage, we anticipate developing several new methods in addition to those described to make the best use of the primary results. We are particularly interested in addressing the problem of ambiguity, where mutually incompatible interpretations are in competition for pieces of observed evidence. Unlike previous levels of analysis, the secondary grouping stage permits the explicit

recognition and resolution of such conflicts. There are also opportunities to add processing on top of this layer, for instance to detect still larger-scale structures, perhaps with some kind of learning scheme.

In conclusion we find that the sinusoid track representation of sound provided for a very natural and straightforward implementation of grouping rules based upon the cues to source formation known from psychoacoustics. Having the results of these first-layer rules allowed us to experiment with secondary grouping methods, which led to reasonably high-level, robust structures identified in real sound examples. We are very hopeful that a refined and expanded set of rules will allow us to build a truly useful model of human auditory event formation.

Acknowledgments

This work was generously supported by the MIT Media Laboratory. The author is in the United States as a Harkness Fellow of the Commonwealth Fund of New York, whose support is gratefully acknowledged.

References

- [1] AS Bregman (1990) *Auditory Scene Analysis*, MIT Press
- [2] MP Cooke (1991) "Modeling auditory processing and organisation," PhD thesis, CS dept, Univ. of Sheffield
- [3] GJ Brown (1992) "Computational auditory scene analysis: A representational approach," PhD thesis CS-92-22, CS dept, Univ. of Sheffield
- [4] HG Okuno, T Nakatani, T Kawabata (1994) "Auditory stream segregation in auditory scene analysis with multi-agent system," in Proc. AAAI-94, Seattle
- [5] DK Mellinger (1991) "Event formation and separation in musical sound," PhD thesis, CCRMA, Stanford Univ.
- [6] K Kashino, H Tanaka (1993) "A sound source separation system with the ability of automatic tone modeling," Proc. Int. Comp. Music Conf, Tokyo
- [7] M Minsky (1986) *The Society of Mind*, Simon and Schuster
- [8] SH Nawab, V Lesser (1992) "Integrated signal processing and understanding of signals," in *Symbolic and knowledge-based signal processing*, ed. AV Oppenheim & SH Nawab, Prentice Hall
- [9] DPW Ellis (1992) "A perceptual representation of audio," MS thesis, MIT EE dept.
- [10] RJ McAulay, TF Quatieri (1986) "Speech analysis/synthesis based on a sinusoidal representation" IEEE Tr. ASSP 34