

Towards a Hardware Realization of Time-Frequency Source Separation of Speech

Naomi Harte* Niall Hurley† Conor Fearon† Scott Rickard†

Abstract — This paper presents preliminary work on a hardware implementation of a source separation algorithm employing time-frequency masking methods. DUET (Degenerate Unmixing Estimation Technique) has previously been shown to achieve excellent source separation in real time in software. The current work is a move towards a hardware realization of DUET that will allow integration of the algorithm into consumer devices. Initial stages involve investigating the performance of DUET when implemented in fixed-point arithmetic and a consideration of algorithmic changes to make DUET more amenable to implementation on a DSP processor. Performance is compared for floating-point and fixed-point implementations. A Weighted K-means clustering algorithm is presented as an alternative to gradient descent methods for peak tracking and demonstrated to achieve excellent performance without adversely affecting computational load. Preliminary performance figures are given for an implementation on a TMS320VC5510 DSK.

1 INTRODUCTION

Blind source separation is a classic problem involving the separation of unknown sources from a number of mixtures. The DUET (Degenerate Unmixing Estimation Technique) algorithm has been widely reported in the literature as a computationally inexpensive method of separating an arbitrary number of sources from just two mixtures [3]. The original algorithm was not done in real-time and involved a two-pass approach where the mixtures were first used to build a histogram in amplitude-delay space [1] and then separated in a second pass through the data. Subsequent work [2] extended the algorithm to real-time operation where data could be processed frame-by-frame. In this case, a gradient descent method was used to track peak delay and amplitude values over time and allow real-time time-frequency masking of the data to yield the original sources.

This current work is motivated by an interest in implementing DUET in real-time on hardware to demonstrate that the algorithm is suitable for integration into low-cost consumer devices, e.g. PDA. To that end, the algorithm has been migrated to a fixed-point implementation and changes to the

algorithm making it more amenable to hardware realization have also been considered. Details of an initial implementation on a TMS320VC5510 DSK are reported.

Section 2 gives a brief overview of the DUET algorithm for those unfamiliar with its operation. Readers will be directed to the relevant literature for greater detail. A system overview is given of how the algorithm operates in real-time in section 3.1. Section 3.2 will outline issues encountered in migrating DUET to fixed-point operation. K-means clustering is then considered as an alternative method for peak tracking in the DUET algorithm. Experiments are reported in section 4 to compare the performance of the floating point and fixed-point algorithm. Results also demonstrate the performance of Weighted K-means clustering for peak tracking. Details of a preliminary porting of the algorithm to a TMS320VC5510 DSK are also given.

2 THE DUET ALGORITHM

2.1 Problem Definition

DUET is a blind source separation technique capable of the separation of N sources from 2 mixtures. The N sources can be defined as $s_1(t), s_2(t) \dots s_N(t)$. Let $x_1(t)$ and $x_2(t)$ be the mixtures derived from these sources as:

$$x_1(t) = \sum_{j=1}^N s_j(t), \quad (1)$$

$$x_2(t) = \sum_{j=1}^N a_j s_j(t - \delta_j), \quad (2)$$

where δ_j is the arrival delay between sensors resulting from the angle of arrival, and a_j is a relative attenuation factor corresponding to the ratio of the attenuations of the paths between sources and sensors. Yilmaz and Rickard [3] have previously presented the theory in detail. However, the result of interest in this paper is the fact that by presuming anechoic conditions and that the source signals are approximately w-disjoint orthogonal, only one source is active at any time-frequency point. By using parameter estimation techniques to estimate delay and attenuation values, it is then possible to

*Dept. of Electrical and Computer Engineering, McMaster University, Ontario, Canada

†Dept. of Electronic and Electrical Engineering, University College, Dublin, Ireland

construct a time-frequency mask M_j that will isolate source j from the mixtures.

2.2 Gradient Descent Search

As outlined by Rickard et al. [2], in order to achieve real-time operation of DUET, a gradient search technique is used for mixing parameter estimation over time. Given initial estimates of the delay and attenuation parameters a cost function $J(\tau)$ can be derived as

$$J(\tau) = \min_{a_1, \delta_1, \dots, a_N, \delta_N} \sum_{\omega} -\frac{1}{\lambda} \ln(e^{-\lambda \rho_1} + \dots + e^{-\lambda \rho_N}), \quad (3)$$

where

$$\rho(a_j, \delta_j, \omega, \tau) \doteq \frac{1}{1 + a_j^2} |X_1(\omega, \tau) a_j e^{-i\omega \delta} - X_2(\omega, \tau)|^2. \quad (4)$$

Given that the number of sources being searched for is known, it is possible to derive updates for the amplitude and delay values $(a_j[k], \delta_j[k])$ from the current frame $\tau_k = k\tau_{\Delta}$ as:

$$a_j[k] = a_j[k-1] - \beta \alpha_j[k] \frac{\partial J(\tau_k)}{\partial a_j}, \quad (5)$$

$$\delta_j[k] = \delta_j[k-1] - \beta \alpha_j[k] \frac{\partial J(\tau_k)}{\partial \delta_j}, \quad (6)$$

where β is a learning rate constant and $\alpha_j[k]$ is a time and mixing parameter dependent learning rate for time index k for estimate j .

3 REAL-TIME OPERATION

3.1 System Overview

The diagram in Figure 1 shows a block diagram of the system used to implement DUET in real time.

With data input at 8KHz, data is processed in frames of 512 samples with a 75% overlap of successive frames. After windowing with a hamming window, the frame is transformed to the frequency domain via a 512-point FFT. Instantaneous delay and amplitude values for this frame are calculated for each time-frequency point as outlined by Rickard et al.[3]. Updated parameter estimates for the actual amplitude and delay values are obtained using the gradient update as explained in equations [5, 6] above.

The masking operation involves calculating, for each time-frequency point, which of the N peak amplitude and delay values each point is closest too. This is done using a simple Euclidean distance measure. In this way each time-frequency point is only assigned to a single source. This “winner take all”

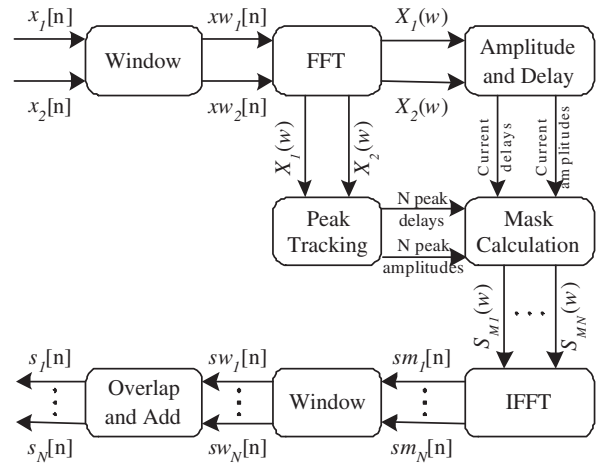


Figure 1: DUET System - Block Diagram

scenario greatly reduces computational complexity and, as reported previously in [2], has little perceivable impact on performance. The N masks are used to derive a time frequency representation of each of the sources as

$$S_{M_j}(\omega, \tau) = M_j(\omega, \tau) X_1(\omega, \tau). \quad (7)$$

An inverse FFT on each of the N masked signals, followed by windowing and overlap-and-add yields a new frame of each of the N source signals.

3.2 Fixed-Point Migration

The aim of this work is to demonstrate the feasibility of incorporating the DUET algorithm into high-end consumer devices. Considering costs, it was hence appropriate that the algorithm be targeted at a fixed-point rather than floating-point processor. The TI C5510 family was chosen as the target processor and the TMS320VC5510 DSK represented a low-cost development platform for developing the algorithm. This device is a 16-bit processor, operating at 200MHz and capable of delivering up to 400 MIPs. The chip has 160K 16-Bit On-Chip RAM and a dual MAC. Full details of the chip and DSK are available online from the at the TI website [4] and Spectrum Digital homepage [5].

A fixed-point implementation equivalent to the floating-point system was carried out in C, as an intermediate step to allow full system testing of the fixed-point migration. Within this system a number of simplifications were made to speed up development. Any functional blocks such as trigonometric functions, FFT, log were not written as full fixed-point libraries as these would be integrated when targeting the board. For system evaluation,

inputs and outputs of these functions were appropriately quantized. This gives a slightly more advantageous performance than the final system but was considered appropriate for development purposes. All other values were stored as 16-bits. For the port to the TMS320VC5510 DSK, the free signal processing libraries supplied by TI were used.

3.3 K-Means Clustering

Initial evaluation of the computational load in DUET revealed that the main part of the effort lay in the evaluation of gradients used in the parameter updates as described in equations [5,6]. Furthermore, as the results section will demonstrate, problems were encountered with the performance of the fixed-point gradient descent. Hence, it was considered worthwhile to investigate alternative methods of tracking the peak values.

The K-means algorithm is a classic technique employed in data clustering problems [6]. The algorithm efficiently partitions the points of a data matrix into K clusters. It achieves this by minimizing (in a least mean squares sense) the sum of distances from each point to its nearest cluster center. Each iteration starts by reassigning points to their nearest cluster center. Each cluster center is then recalculated as the mean of all points which have been assigned to it. This process is repeated until the cluster centers converge according to the chosen criterion.

The use of a histogram space has proved to be very powerful in DUET [1]. Rather than searching for peaks in the entire amplitude-delay space, the data is placed into a bounded histogram with a finite number of bins. The limitation of this method for real-time operation is the use of a two-pass approach. In the current work, K-means clustering is used to allow peak tracking in histogram space in real-time. A weighted version of the K-means algorithm is performed on the histogram. For each frame of data, the (weighted) histogram is updated with the powers of the corresponding time-frequency points. The histogram bin centers are passed to K-means and each point is weighted by the height of that histogram bin. In this way, peak-tracking updates are calculated using information from all previous frames. As the results section demonstrates, this method yields very accurate peak estimates.

4 EXPERIMENTS

This section outlines a number of experiments designed to test the fixed-point performance of DUET. Performance is compared to the original

floating-point algorithm. Results on the use of the Weighted K-means algorithm for peak tracking are also given.

4.1 Fixed-Point Performance of DUET

The difference in peak estimates for amplitude and delay were compared for the floating-point and fixed-point implementations of DUET. The table below details the average percentage error in the fixed-point estimate, referenced to the floating-point value obtained for a mixture of two sources.

Peak	% Error
Amplitude	11.15
Delay	94.62

Table 1: Percentage Error in Delay and Amplitude Estimates for Fixed-Point System

A significant error has been introduced, particularly in the delay estimate. Examining the evolution of the delay estimate over time in Figure 2, it is clear that the delay is not converging to the true value. Closer investigation has shown that this error is largely attributable to underflow in the derivative values for the gradient update. This arises due to the large number of successive multiplies used. Even with appropriate scaling, a significant number of derivative values simply tend to zero. Clearly this method of peak tracking is problematic using fixed-point arithmetic. It is anticipated that the use of weighted K-means will alleviate this issue.

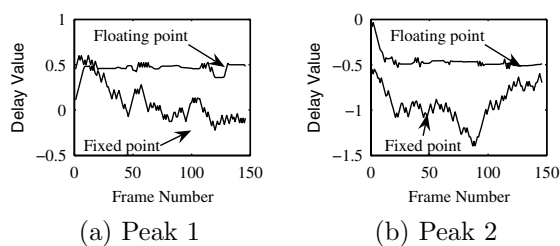


Figure 2: Emerging Delay Estimates for 2 Sources.

An initial port to the TMS320VC5510 DSK has yielded an upper estimate of 18 MIPS for the DUET functionality. It should be noted that this is an unoptimized port which incorporates the gradient descent peak tracking (including estimates of divide functionality) and current indications suggest a final figure of 5 MIPS as highly achievable. Work is ongoing on optimizing the performance of the algorithm in hardware. This involves migrating

to proprietary libraries for trigonometric functions, the FFT and incorporating the weighted K-means method of peak tracking.

4.2 Performance of Weighted K-Means

As an indicator of the comparative complexity of each gradient descent and Weighted K-means for peak tracking, the number of adds and multiplies for a speech file of length N samples is shown in Table [2] for both algorithms.

	ADDS	MULTIPLIES
Weighted K-Means (128^2 bins)	$6671N$	$4533N$
Weighted K-Means (64^2 bins)	$1555N$	$1062N$
Weighted K-Means (32^2 bins)	$318N$	$227N$
Weighted K-Means (16^2 bins)	$69N$	$53N$
Gradient Descent	$44N$	$165N$

Table 2: No. of adds and multiplies for Weighted K-means and Gradient Descent algorithms.

Clearly, Weighted K-means is highly dependent on the number of bins used in the histogram space. Reducing the number of histogram bins by a power of two reduces the number of adds and multiplies by the same factor. For the separation of only two sources, it has been found that smaller histogram spaces of 16×16 bins still yield good performance. However, a histogram space of 128×128 would be required for the case of more than 4 sources. An important advantage of Weighted K-means is that it is possible to completely eliminate the need for any divides. This is not possible with the current formulation of the gradient descent estimate updates.

The error in peak estimates for a two-source mixture at each frame is shown in Figure 3. The Weighted K-means with 128^2 and 16^2 histogram bins are represented by the solid and dotted lines respectively. The Gradient Descent algorithm is represented by the dash-dotted line. It is clear that the peak estimates produced by Weighted K-means are considerably closer to the true values and that this small error is maintained even at coarser resolution in the histogram space when only two sources are present.

5 CONCLUSIONS

This paper has presented initial work on the migration of the DUET algorithm to a fixed-point implementation in hardware. Significant problems were encountered in migrating to a fixed-point implementation of DUET incorporating Gradient Descent peak tracking. Weighted K-means clustering is shown to outperform Gradient Descent for

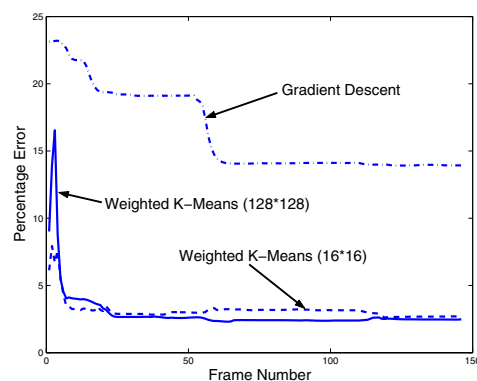


Figure 3: Total percentage error plots for Weighted K-means and Gradient Descent.

amplitude and delay peak tracking without significant adverse effects on computational load. Work is ongoing on the fixed-point implementation to integrate Weighted K-means clustering and to optimize performance on the DSP. Other algorithmic enhancements currently being considered include the exploitation of the properties of speech and improvements in performance in echoic conditions.

Acknowledgments

Work funded by the Enterprise Ireland Proof of Concept fund: Contract PC/2004/347.

References

- [1] A. Jourjine, S. Rickard, and O. Yilmaz, "Blind Separation of Disjoint Orthogonal Signals: Demixing N Sources from 2 Mixtures", IEEE conference on Acoustics, Speech, and Signal Processing (ICASSP2000), Vol 5, pp 2985–2988, Istanbul, Turkey, June 2000.
- [2] S. Rickard, R. Balan, and J. Rosca, "Real-Time Time-Frequency Based Blind Source Separation", 3rd International Conference on Independent Component Analysis and Blind Source Separation (ICA2001), San Diego, CA, December 9-12, 2001
- [3] O. Yilmaz and S. Rickard, "Blind Separation of Speech Mixtures via Time-Frequency Masking", IEEE Transactions on Signal Processing, Vol. 52(7), pp 1830–1847, July 2004.
- [4] "TMS320C5000(tm) Platform Overview", <http://dspvillage.ti.com/>
- [5] "DSP Starter Kit for the TMS320VC5510", <http://www.spectrumdigital.com/>
- [6] G.A.F. Seber, "Multivariate Observations", Wiley, New York, 1984.