

# An Automatic Kurtosis-Based *P*- and *S*-Phase Picker Designed for Local Seismic Networks

by Christian Baillard, Wayne C. Crawford, Valérie Ballu, Clément Hibert, and Anne Mangeney

**Abstract** We present an automatic *P*- and *S*-wave onset-picking algorithm, using kurtosis-derived characteristic functions (CF) and eigenvalue decompositions on three-component seismic data. We modified the kurtosis CF (Saragiotis *et al.*, 2002) to improve pick precision by computing the CF over several frequency bandwidths, window sizes, and smoothing parameters. Once phases are picked, our algorithm determines the onset type (*P* or *S*) using polarization parameters, removes bad picks using a clustering procedure and the signal-to-noise ratio (SNR) and assigns a pick quality index based on the SNR.

We tested our algorithm on data from two different networks: (1) a 30-station,  $100 \times 100$  km array of mostly onland wideband seismometers in a subduction context and (2) a four-station,  $7 \times 4$  km array of ocean-bottom seismometers over a midocean ridge volcano. We compared picks from the automatic algorithm with manual and short-term average/long-term average (STA/LTA)-based automatic picks on subsets of each dataset. For the larger array, the automatic algorithm resulted in more locations than manual picking (133 versus 93 locations out of 163 total events detected), picking as many *P* onsets and twice as many *S* onsets as with manual picking or the STA/LTA algorithm. The difference between manual and automatic pick times for *P*-wave onsets was  $0.01 \pm 0.08$  s overall, compared with  $-0.18 \pm 0.19$  s using the STA/LTA picker. For *S*-wave onsets, the difference was  $-0.09 \pm 0.23$  s, which is comparable to the STA/LTA picker, but our picker provided nearly twice as many picks. The pick accuracy was constant over the range of event magnitudes ( $0.7$ – $3.7 M_l$ ). For the smaller array, the time difference between our algorithm and manual picks is  $0.04 \pm 0.17$  s for *P* waves and  $0.07 \pm 0.08$  s for *S* waves. Misfit between the automatic and manual picks is significantly lower using our procedure than using the STA/LTA algorithm.

## Introduction

Earthquake hypocenter locations are needed to map existing faults and to document their activity, both of which are of prime importance in defining hazards and forecasting events. The number of seismic stations around the world is rapidly growing, which should provide much more detailed information about seismically active regions, but only if seismic events can be accurately and uniformly picked on each station. Automatic picking procedures (APP) are needed to handle the larger datasets and they must be reliable, precise, and capable of distinguishing different phase onsets and adapting to different site and/or instrument characteristics. Compared with manual picking, APP save time and should be more consistent, because manual picks can differ between analysts (Freedman, 1966; Zeiler and Velasco, 2009).

In seismology, the most commonly used event detection algorithm is the short-term average/long-term average (STA/LTA) detector proposed by Allen (1982), which is based on the ratio of the two averages calculated on sliding

windows over the trace. This algorithm is rapid and remains useful for detecting events in continuous databases, but it generally gives significantly different results from manual picking (Saragiotis *et al.*, 2002). The STA/LTA algorithm can be applied to raw traces or to derived traces function called characteristic functions (CF). Baer and Kradolfer (1987) improved the STA/LTA by introducing the envelope function as the CF and by using a dynamic threshold to detect signals buried in noise.

Takanami and Kitigawa (1993), Sleeman and van Eck (1999), and Leonard and Kennett (1999) proposed another approach for automatic picking, derived from auto-regressive (AR) methods. These methods involve calculating AR models for two stationary segments. These two models will be most different, when one contains only seismic noise and the others mostly signal. The Akaike Information Criterion (AIC; Akaike, 1974), which indicates the unreliability of the model fit, is then used to precisely pick the onset. These

methods have not been applied to many datasets, because the calculations are computationally expensive and because some arrivals are not necessarily associated with important changes in their frequency content and so are not detected.

Saragiotis *et al.* (2002) was one of the first to apply higher order statistical functions to seismic traces, introducing the skewness and the kurtosis functions to phase picking. Küperkoch *et al.* (2010) added the AIC to Saragiotis' method and developed a quality-weighting scheme for picks. Both methods provide better accuracy than the STA/LTA and the Baer and Kradolfer (1987) methods.

Nippres *et al.* (2010) introduced the concept of tandem automatic pickers by combining either the STA/LTA or the predominant period time-domain estimation method ( $T^{\text{pd}}$ ) (Hildyard *et al.*, 2008) with the kurtosis characteristic function developed by Saragiotis. Improvement of picking accuracy is significant, but a parameter optimization step is required for each station prior to automatic picking.

We present a new algorithm for automatically picking onsets based on the kurtosis method. We focus on picking accuracy and on the simplicity of implementation, by eliminating operator-intensive phases. We pay particular attention to automatically identifying both  $S$  and  $P$  waves, because combined  $S$  and  $P$  picks can significantly improve hypocenter locations (Gomberg *et al.*, 1990). The steps of the method are (1) computation of the kurtosis CF for each component, (2) modification of the CF to improve pick accuracy, (3) differentiation between  $S$  and  $P$  waves using the wavefield polarization, (4) estimation of the quality index using the signal-to-noise ratio (SNR), (5) rejection of erroneous picks using clustering and distribution analysis, and (6) calculation of the signal amplitude. Each step has only a few variable parameters, simplifying the algorithm's use.

The algorithm relies on three-component data in stage 3, to differentiate  $P$  and  $S$  waves and to reject surface waves and noise spikes. Many present-day networks use three-component seismic sensors, but we also show a case in which we modified the algorithm to identify  $P$  and  $S$  waves using single component data.

We applied our scheme to a network of 30 wideband seismometers in the Vanuatu region, where a large number of earthquakes are generated by the subduction of the Australian Plate under the North Fiji basin (Pelletier *et al.*, 1998). We also applied the automatic picker to a network of four short-period ocean-bottom seismometers (OBS) deployed in the Azores region to study micro seismicity in the vicinity of the mid-Atlantic ridge (MAR) in the framework of the monitoring of the mid-Atlantic ridge (MoMAR) project (Colaço *et al.*, 2011; Crawford *et al.*, 2013). In both cases, we compared automatic and manually picked onsets to test the quality of the automatic picker and the efficiency of the quality index assignment, and we compared the performance of our picker with that of an STA/LTA picker.

## Mathematical Background

### Kurtosis

The kurtosis is a statistical value characterizing the shape of a given distribution. It is a positive scalar defined as the standardized fourth moment about the mean. Using a probabilistic notation, the kurtosis  $K$  is

$$K \equiv \frac{E[(X - \mu)^4]}{\{E[(X - \mu)^2]\}^2} = \frac{m_4}{\sigma^4}, \quad (1)$$

in which  $X$  is a random variable,  $E$  is the expectation operator,  $\mu$  is the mean,  $m_4$  is the fourth central moment, and  $\sigma$  is the standard deviation. When considering a numerical signal of  $n$  samples, represented as  $x = \{x_1, x_2, \dots, x_n\}$ , the discretized form of equation (1) is

$$K = \frac{\frac{1}{n} \sum_{i=1}^{n+1} (x_i - \bar{x})^4}{\left[\frac{1}{n} \sum_{i=1}^{n+1} (x_i - \bar{x})^2\right]^2}, \quad (2)$$

in which  $\bar{x}$  is the mean over the  $n$  samples.

The kurtosis is 3 for a normal (Gaussian) distribution (DeCarlo, 1997) and generally increases for a non-Gaussian distribution. Seismic-wave onsets temporarily generate a non-Gaussian wavefield that rapidly increases the kurtosis, which we can use to accurately pick the onset times.

### Polarization Analysis

The polarization of wave onsets depends on the wave type (e.g., surface or body waves) and orientation. We describe here two polarization parameters that can be used to distinguish between  $P$  and  $S$  waves, surface waves, and noise. Three-component data are required to calculate these parameters: we assume the data are composed of three orthogonal ground-motion recordings corresponding to the east, north, and vertical components (respectively noted  $X$ ,  $Y$ , and  $Z$ ).

The first step is to compute the  $3 \times 3$  covariance matrix over the three components using an  $N$ -samples sliding window.

$$\mathbf{C} = \begin{pmatrix} \text{Cov}(X, X) & \text{Cov}(X, Y) & \text{Cov}(X, Z) \\ \text{Cov}(Y, X) & \text{Cov}(Y, Y) & \text{Cov}(Y, Z) \\ \text{Cov}(Z, X) & \text{Cov}(Z, Y) & \text{Cov}(Z, Z) \end{pmatrix}, \quad (3)$$

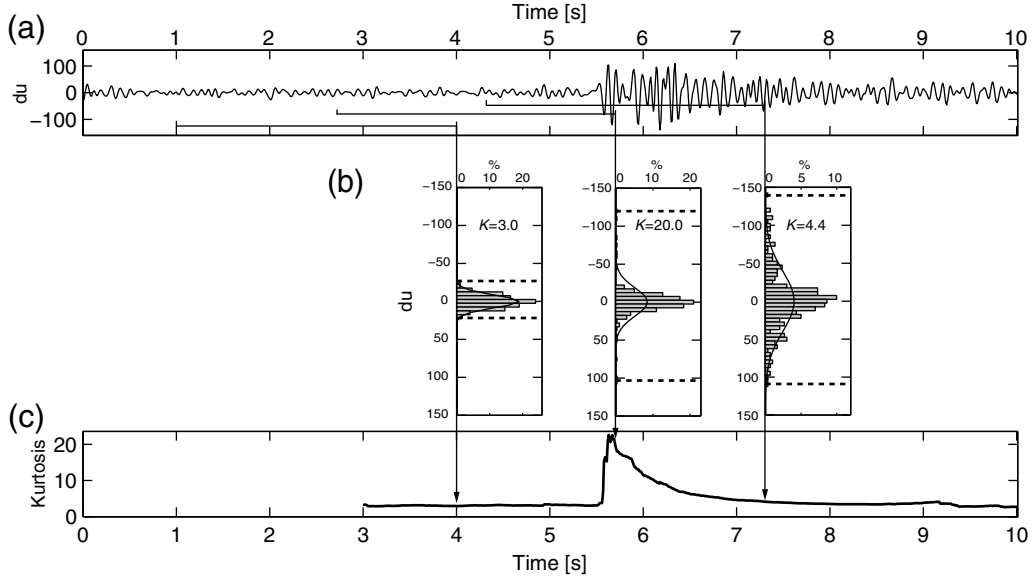
in which the covariance between  $X$  and  $Y$  is defined as

$$\text{Cov}(X, Y) = \frac{1}{N} \sum_{i=1}^N x_i y_i. \quad (4)$$

The next step is the principal component analysis using eigenvalue decomposition. We search for the eigenvalues ( $\lambda_1, \lambda_2, \lambda_3$ ) associated with the unit eigenvectors ( $U_1, U_2, U_3$ ) that satisfy

$$\mathbf{C} \times U_i \equiv \lambda_i U_i \quad i \in [1, 2, 3]. \quad (5)$$

Because the covariance matrix in equation (3) is symmetric and composed of real elements, the eigenvalues are



**Figure 1.** Example of a kurtosis characteristic function calculated on a seismic trace. In this example, a 5 s window is used. (a) Seismic trace (filtered between 3 and 45 Hz for clarity). The horizontal bars indicate the three kurtosis windows examined below. (b) Kurtosis and histograms of sample distributions for the three kurtosis windows. The solid line shows the best-fitting Gaussian curve, and dashed lines show the outer bounds of sample values. (c) The resulting kurtosis characteristic function.

real and the eigenvectors form an orthogonal base. We organize eigenvalues so that  $\lambda_1 \geq \lambda_2 \geq \lambda_3$ .  $\lambda_i U_i$  are called the principal axes of the polarization ellipsoid.

We can now calculate two polarization parameters: the degree of rectilinearity (Jurkevics, 1988) and the dip of maximum polarization (Vidale, 1986). The first polarization parameter, the degree of rectilinearity, is defined as

$$\text{Rec} \equiv 1 - \left( \frac{\lambda_2 + \lambda_3}{2\lambda_1} \right). \quad (6)$$

For example,  $\text{Rec} = 0$  for circular polarization ( $\lambda_1 = \lambda_2 = \lambda_3$ ), and  $\text{Rec} = 1$  for rectilinear polarization ( $\lambda_1 = 1$  and  $\lambda_2 = \lambda_3 = 0$ ).  $P$  and  $S$  waves are body waves, so their degree of rectilinearity is close to 1.

The second polarization parameter, the dip of maximum polarization, is defined as

$$\text{Dip} \equiv \tan^{-1} \left( \frac{U_1(3)}{\sqrt{U_1(2)^2 + U_1(1)^2}} \right). \quad (7)$$

Possible values range from  $-90^\circ$  to  $+90^\circ$ . A horizontal maximum polarization vector has dip =  $0^\circ$ .

### The Characteristic Function

We describe here the steps we follow to create the final CF used in the APP. The goal of this CF is to allow accurate and precise picking of all onsets on a section of the trace. We consider a single seismic trace represented by  $x = \{x_1, x_2, \dots, x_n\}$ . The CF is calculated over a sliding window on the signal: let  $T$  be the size, in seconds, of this window. The number of samples in the windows is therefore  $N =$

$(T/\Delta t) + 1$ , in which  $\Delta t$  is the sample interval. The central moment of order  $d$  at sample  $k$  can be written as (Küperkoch et al., 2010)

$$m_d(k) = \frac{1}{N} \sum_{i=1}^N (x_{k-i+1} - \bar{x}_k)^d \quad \text{with } k \in [1, \dots, n], \quad (8)$$

in which  $\bar{x}_k$  represents the mean of the signal from sample  $(k - N + 1)$  to  $k$ . We see here the dependency of the central moment on the size of the window  $N$ : the bigger  $N$  is, the less sensitive the central moment is to transient variations within the trace.

To reduce computation time, we transform the signal into a zero-mean process, which reduces equation (9) to

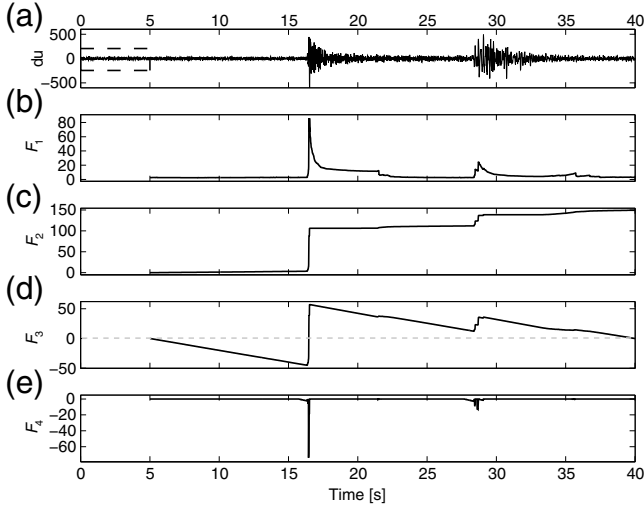
$$m_d(k) \approx \frac{1}{N} \sum_{i=1}^N (x_{k-i+1})^d = \bar{x}_d(k). \quad (9)$$

We then define the first kurtosis CF at sample  $k$  as

$$F_1(k) = \frac{\bar{x}_4(k)}{[\bar{x}_2(k)]^2}. \quad (10)$$

When the sliding window integrates the first samples of a seismic phase onset, the change in the signal distribution from Gaussian to non-Gaussian causes  $F_1$  to increase (Fig. 1).

The best alignment with manual picks is obtained when the automatic pick corresponds to the very beginning of this increase, before its maximum value. Because no simple mathematical tools exist to directly detect this change of behavior, we apply a succession of transformations to the



**Figure 2.** The stages involved in calculating the characteristic function for picking onsets: (a) filtered zero-mean  $Z$  seismogram displaying two strong onsets; (b)  $F_1$ , kurtosis; (c)  $F_2$ , removal of all negative slopes; (d)  $F_3$ , linear correction; and (e)  $F_4$ , pushing down all values by the amplitude of the following maximum, then removing all positive values.

initial CF in order to isolate the initial onset and to identify the strongest onsets (Fig. 2).

The first transformation essentially cleans the initial CF of all strictly negative gradients (Fig. 2c), because only positive gradients characterize the transition from noise to a coherent signal. The transformation is

$$F_2(k+1) = F_2(k) + \delta(k) \times dF_1(k)$$

$$\text{with } \begin{cases} F_2(1) = F_1(1) \\ dF_1(i) = F_1(i+1) - F_1(i) \\ \delta(i) = 1 & \text{if } dF_1(i) \geq 0 \\ \delta(i) = 0 & \text{else} \end{cases} \quad (11)$$

The second transformation removes a linear trend from  $F_2$ , so that the first and last values equal zero (Fig. 2d). In this way, the onsets become local minima of the CF. The transformation is

$$F_3(k) = F_2(k) - [a \cdot (k-1) + b]$$

$$\text{with } \begin{cases} a = \frac{F_2(n) - F_2(1)}{n-1} \\ b = F_2(1) \end{cases} \quad (12)$$

The final transformation makes the amplitude of the minima amplitude scale with the total change in the kurtosis that follows, so that the greatest minima correspond to the greatest onset strengths (Fig. 2e). The transformation pushes the CF values down by the amplitude of the next maximum and sets remaining positive values to zero:

$$T(k) = F_3(k) - M_{i+1} \quad \text{if } k \in [s_i, s_{i+1}] \quad (13)$$

and

$$F_4(k) = T(k) \quad \text{if } T(k) < 0, 0 \text{ otherwise,} \quad (14)$$

in which  $\{M_1, M_2, \dots, M_m\}$  are the local maxima of  $F_3$ , located at samples  $\{s_1, s_2, \dots, s_m\}$ .

Picking minima on  $F_4$  gives a good first estimate of phase onsets, but the picker accuracy is considerably improved by adding two intermediate stages: (1) averaging of  $F_3$  over multiple window lengths and frequency bandwidths; and (2) sequential onset picking—from long-to-short time scales, using smoothing windows on  $F_4$ .

The first stage addresses a common problem of automatic pickers: their accuracy and reliability depend on the frequency bandwidth used to filter the data and the size of the sliding window used to compute the CF. We do not know *a priori* the frequency of the event onsets, and this frequency may change between events. Using only one frequency-window pair, we may lose important information in the signal and introduce frequency-dependent artifacts. To avoid this problem, we compute  $F_3$  over  $p$  different frequency bands (parameter  $BW$ , Tables 1 and 2) and  $l$  different sliding windows (parameter  $WS$ , Tables 1 and 2) and then calculate the average function ( $F'_3$ ) over the  $l \times p$  resulting CFs. This average function reduces artifacts associated with the individual window lengths and frequency bands and gives a much clearer pick of the phase onsets (Fig. 3).

The second stage helps to identify the first arrival in complicated or emergent onsets. In these cases, an onset creates several closely grouped minima on  $F_4$  rather than one big minimum. This is particularly common for emergent arrivals, whose kurtosis increases progressively before reaching its maximum value. To accurately identify the first minimum in the group we adopt a long-to-short time-scale approach. We apply several smoothing windows to  $F'_3$ , then compute  $F_4$  for each smoothed function. The set of smoothing windows is user defined, but the longest window should not exceed the minimum separation between  $P$  and  $S$  onsets. We detect all onsets on the smoothest  $F_4$ , then pick the closest onsets on the next smoothest version of  $F_4$ , and repeat until we pick the corresponding onsets on the least smoothed version of  $F_4$ . Figure 4 shows functions  $F'_3$  and  $F_4$  for different smoothing parameters when zooming on an onset. This stage significantly improves the picking reliability and allows our algorithm to be applied to a wide range of non-impulsive onsets.

### $P$ – $S$ Characterization

Once we have created the CF that best allows us to pick the onsets on a single trace, we need to determine whether those onsets correspond to surface waves or body waves or if they are not seismic waves at all. If the onsets correspond to body waves, we must differentiate  $P$  and  $S$  onsets. We assume the rays arrive at the stations with a small incidence angle to the vertical axis, meaning the  $P$  waves will have mostly vertical motions and the  $S$  waves mostly horizontal motions. This assumption will be true on most stations if

Table 1  
Description of Parameters Used in the APP

Step	Parameter	Description	Comments
SNR-based trace selection	$N_b$ (s)	SNR prewindow	Long enough to represent preonset noise energy
	$N_a$ (s)	SNR postwindow	Long enough to cover onset energy
	$W_c$ (s)	New analysis window	Section of the trace that contains the $P$ - $S$ onsets
	$v_c$ (-)	Maximum number of desired onsets	If (number of onsets in $W_c$ ) > $v_c$ then the trace is rejected, usually $\leq 2$
	$T_{\text{SNR}}$ (%)	SNR threshold (% of max SNR)	20%, normally do not change
CF construction	$\text{SNR}_{\text{min}}$ (-)	Minimum SNR to accept pick	$\geq 1$ , generally less than 5
	$BW$ (Hz)	$F_3$ Frequency band	Usually two: 1) full useful data bandwidth and 2) bandwidth of typical arrivals
	$WS$ (s)	$F_3$ Window lengths	Several equally spaced up to average $P$ - $S$ delay
	$N_s$ (samples)	Smoothing window lengths	From 1 to just below the average $P$ - $S$ delay divided by the sampling rate, spaced closely enough to not jump onsets between smoothing windows
Polarization analysis	$T_P$ (-)	$P$ threshold	Between 0.5 and 1: the higher the numbers the more selective the analysis
	$T_S$ (-)	$S$ threshold	Approximately $-T_P$
	$W_{\text{pol}}$ (s)	Window length for polarization analysis	4-seconds, normally do not change
	$\alpha$ (-)	DR calculation weighting factor	Between 1 and 2: higher is more selective
Pick quality analysis	$P0; P1; P2; P3$ (-)	$P$ -SNR ratios for pick weight numbers	$10 \leq P0 \leq 15$ (higher is more selective), $P3 = \text{SNR}_{\text{min}} - P1$ and $P2$ evenly spaced between
	$M_f$ (-)	( $S$ -SNR ratios)/( $P$ -SNR ratios)	1.25, normally do not change
Clustering rejection	$W_P$ (s)	$P$ clusters cutoff	Greater than largest typical $\Delta t$ between $P$ -phase arrivals across network
	$W_S$ (s)	$S$ clusters cutoff	Approximately 1.7 times $W_P$
Amplitude calculation	$A_{PS}$ (s)	$P$ - $S$ window length	Typical $P$ - $S$ time in dataset
	$A_0$ (s)	Onset window length	Typical length of energy envelope after onsets

Table 2  
APP Parameter Values for the Vanuatu and MAR Datasets

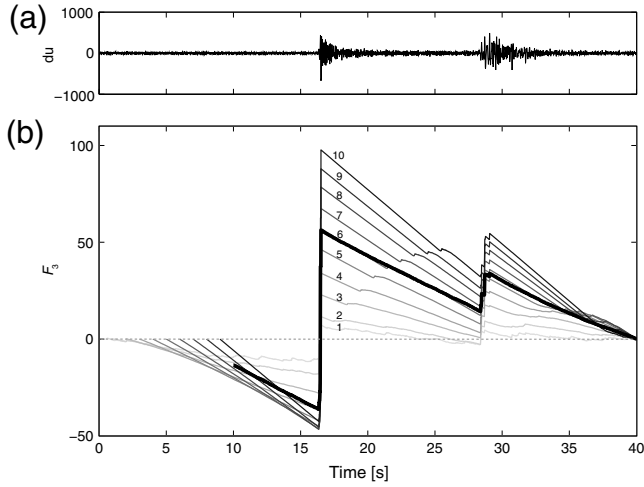
Parameter	Vanuatu	MAR
$N_b$ (s)	2	0.5
$N_a$ (s)	2	0.5
$W_c$ (s)	30	15
$v_c$ (-)	2	2
$T_{\text{SNR}}$ (%)	20	20
$\text{SNR}_{\text{min}}$ (-)	5	3
$BW$ (Hz)	[3–18] and [3–45]	[5–30] and [3–18]
$WS$ (s)	2, 4, and 6	0.5, 1, and 1.2
$N_s$ (samples)	1, 2, 3, 6, 8, 10, 20, 40, 50, 70, and 80	1, 2, 3, 6, 8, 10, 20, 40, and 50
$T_P$ (-)	0.4	/
$T_S$ (-)	-0.4	/
$W_{\text{pol}}$ (s)	4	/
$\alpha$ (-)	1.3	/
$P0; P1; P2; P3$ (-)	(15, 11, 7, $\text{SNR}_{\text{min}}$ )	(11, 8, 5, $\text{SNR}_{\text{min}}$ )
$M_f$ (-)	1.25	1.25
$W_P$ (s)	10	2
$W_S$ (s)	30	3
$A_{PS}$ (s)	10	5
$A_0$ (s)	10	5

the network geometry is adequate and can be verified for each station *a posteriori*, after hypocenter inversion. We compute the dip and rectilinearity using a sliding window of size  $W_{\text{pol}}$  (Tables 1 and 2) on the three components  $X$ ,  $Y$ , and  $Z$ . Further details about how to correctly choose

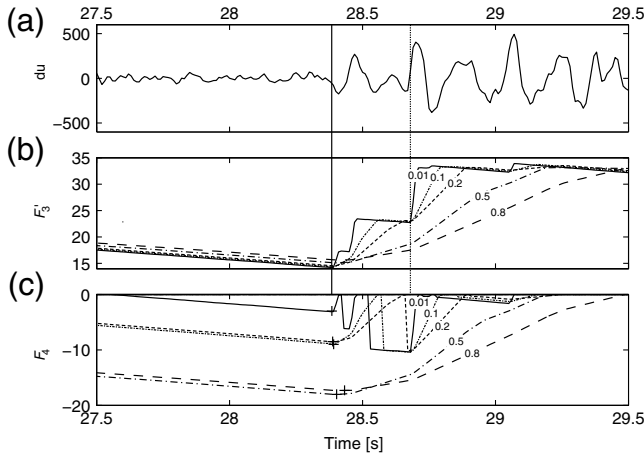
the polarization window size are given at the end of this section. We then define the dip-rectilinearity function as

$$\text{DR}(k) \equiv \text{Rec}(k) \times \text{sign}[\alpha \times \sin(|\text{Dip}(k)|) - \text{Rec}(k)], \quad (15)$$

in which the parameter  $\alpha$  (Tables 1 and 2) is a user-defined weighting factor between 1 and 2 that depends on the clarity of the dip and rectilinearity: a value of 2 would be appropriate for perfectly polarized data, and 1 corresponds to poorly polarized data. The rectilinearity parameter helps to separate body waves ( $\text{Rec}(k)$  close to 1) from background noise and surface waves, whereas the dip parameter distinguishes between  $P$  and  $S$  waves.  $P$  waves have high dip values due to their polarization along the near-vertical incidence ray ( $\sin(|\text{Dip}(k)|)$  is close to 1). In contrast,  $S$  waves are polarized horizontally ( $SH$  and  $SV$ ), and  $\sin(|\text{Dip}(k)|)$  is close to 0. In a perfect case, we could take  $\alpha = 2$ , and  $\text{DR}(k)$  will be 1 for  $P$  waves and -1 for  $S$  waves. Practically, because the waves are not perfectly linearly polarized and the incidence is not exactly vertical, we have to choose a smaller  $\alpha$  and define two thresholds,  $T_P$  and  $T_S$  (Tables 1 and 2), with  $0 < T_P < 1$  and  $T_S \approx -T_P$ . The onset is declared  $P$  if the average of the window centered on the pick is greater than zero and  $\text{DR} > T_P$  somewhere in the window. The onset is declared  $S$  if the average of the window centered on the pick is less than zero and  $\text{DR} < T_S$  somewhere in the window. If neither case is satisfied, the pick is rejected (Fig. 5).

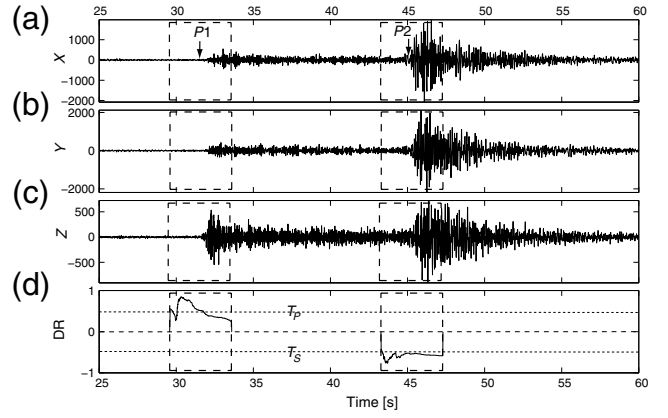


**Figure 3.** Illustration of the multiple window length stage: (a) filtered zero-mean  $Z$  seismogram; and (b)  $F_3$  for each window (thin lines) and the average  $F_3$  (thick line). The numbers next to each thin line indicate the corresponding window length in seconds.



**Figure 4.** Illustration of the long-to-short time-scale approach for identifying first arrivals. (a) Zoom on the second ( $S$  wave) onset of the trace in Figure 3. (b)  $F_3$  computed using five different smoothing windows: 0.01, 0.1, 0.2, 0.5, and 0.8 s. (c)  $F_4$  calculated from each  $F_3$ . Crosses represent the local minima that lead to the final pick (solid line). The dotted line marks the pick that would have been made on the short time-scale data without using this approach. The time difference between the minima picked on the shortest time scale data before and after applying our approach is 0.3 s. The time difference between minima picked on the longest and shortest time scale is 0.03 s.

We performed several tests on the polarization window size to identify the best configuration. Figure 6 shows the influence of the window size on the rectilinearity parameter. The longer the window is, the smoother the rectilinearity function and the easier it is to characterize the waveform type (Fig. 6c,d), but a window longer than the  $P$ - $S$  delay will smear the effects of these two arrival types. To find the optimal value, we calculated the noise level for different window sizes (Fig. 6e). The noise level for a specific window



**Figure 5.** Example of the dip-rectilinearity parameter used to identify seismic phases and reject noise and surface wave picks (Vanuatu dataset). (a-c) Filtered  $X$ ,  $Y$ , and  $Z$  traces showing the 4 s windows used for polarization analysis (centered on arrivals  $P1$  and  $P2$ ). (d) Dip-rectilinearity parameter (DR). Above  $T_P$ , the pick is declared as  $P$ ; below  $T_S$ , the pick is declared as  $S$ . DR extrema do not correspond to wave onsets because of the shifting induced by the polarization analysis.

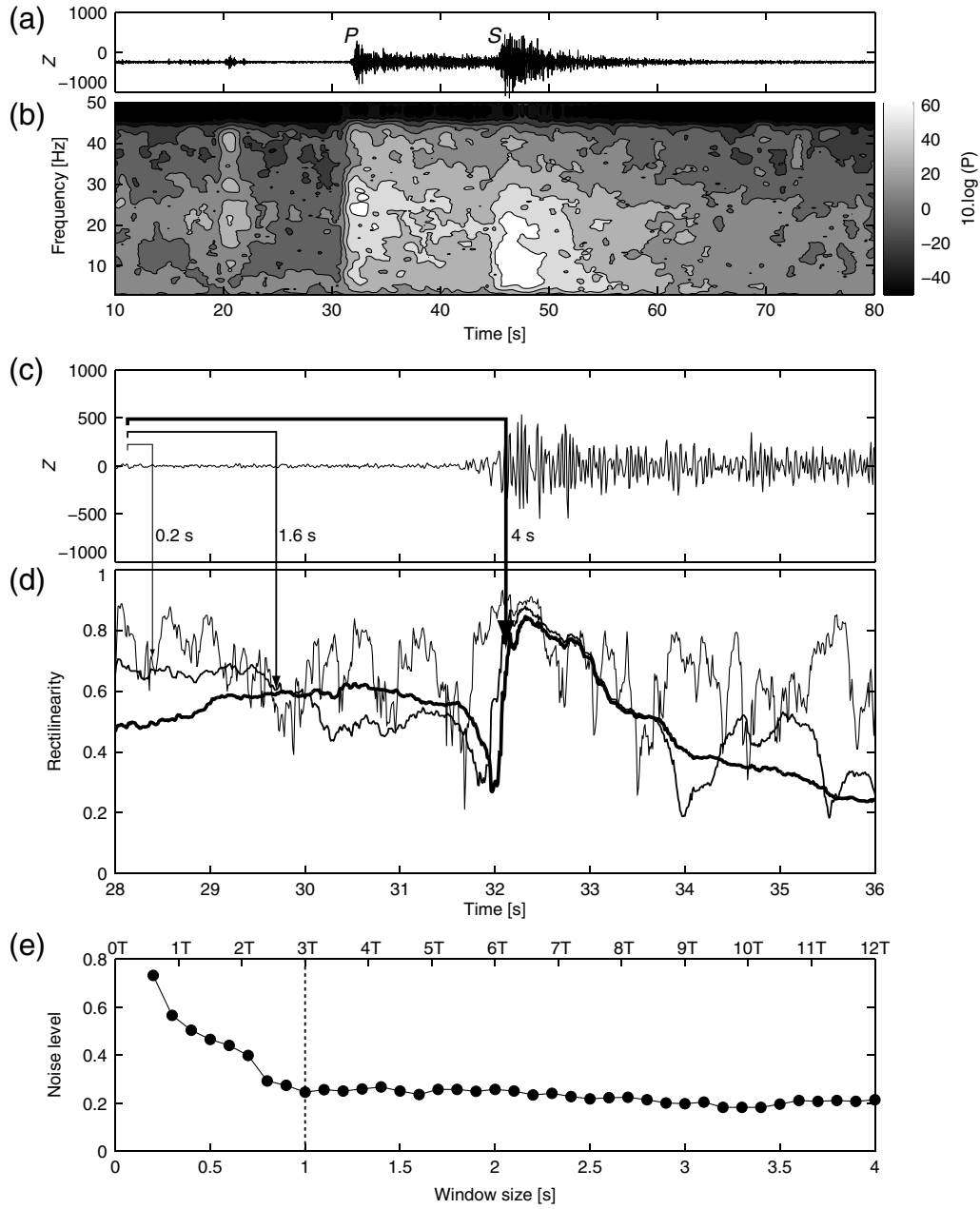
length is obtained by computing what portion of the rectilinearity signal is above the threshold before the arrival onset. This noise level decreases with increasing window sizes until reaching a constant value of 0.2 for windows longer than  $3T$  (Fig. 6e), in which  $T$  is the longest period of energy in the onset (3 Hz for Vanuatu data). For reliable polarization characterizations, we use windows longer than  $3T$  but shorter than the most  $P$ - $S$  delays.

### Automatic Picking Procedure

We describe here the entire event extraction and automatic picking procedure, starting with raw continuous seismic traces and ending with a catalog of picked phase times and amplitudes for each event at each station.

Before running the APP, we extract record sections with a high probability of containing an event from the continuous database. Many standard routines exist for this phase: we used SEISAN's CONDET routine (Havskov and Ottemöller, 2010), which is based on an STA/LTA detection algorithm run on the vertical channel of each station and which extracts a record for each case that has more than a threshold number of stations trigger in a specified time interval. We generally extract a longer window than is strictly necessary, because the first step of the APP will optimize the window length.

The first step of the APP is to reduce the size of the analysis window to an optimal length (parameter  $W_c$ , Tables 1 and 2). This speeds up subsequent processing and helps to avoid mispicks by reducing the chance of having multiple events in one window. The algorithm calculates a simplified version of  $F_4$  (equation 15) for all traces on all stations, picks the biggest onset on each trace, and detects the region of maximum pick density using a clustering analysis.



**Figure 6.** Effect of the polarization window size on the rectilinearity parameter. The seismogram is taken from the Vanuatu dataset. (a) Filtered Z trace containing both P and S waves and (b) the associated spectrogram. (c) Zoom in of the Z trace around the P arrival. (d) Rectilinearity parameter computed for three different window sizes. For the smallest window, the signal is buried in noise. Signals computed with the two longest windows show similar patterns. (e) Noise level of the rectilinearity function of window size, expressed in seconds, and also in number of periods  $T$ ,  $T$  being the period associated with the lowest frequency contained in the onset (3 Hz in this case). The dashed line represents minimum polarization window size over which the rectilinearity is reliable.

All subsequent analysis will be performed on the  $W_c$ -second window centered on this region.

The algorithm now processes each station separately. It first analyzes the onset's SNR using the following expression for pseudoenergy:

$$\text{Energy} = X^2 + Y^2 + Z^2, \quad (16)$$

in which  $X$ ,  $Y$ , and  $Z$  are the three components of the signal. The SNR at sample  $k$  is defined as

$$\text{SNR}(k) \cong 20 \times \log \left( \frac{\text{Energy}_{k \rightarrow k + N_a}}{\text{Energy}_{k - N_b \rightarrow k}} \right). \quad (17)$$

The numerator represents the mean energy from  $k$  to  $k + N_a$ , and the denominator represents the mean energy from  $k - N_b$  to  $k$ .  $N_a$  and  $N_b$  (Tables 1 and 2) are user defined and generally equal:  $N_b$  must be long enough to correctly represent the noise level before the onset, and  $N_a$  must

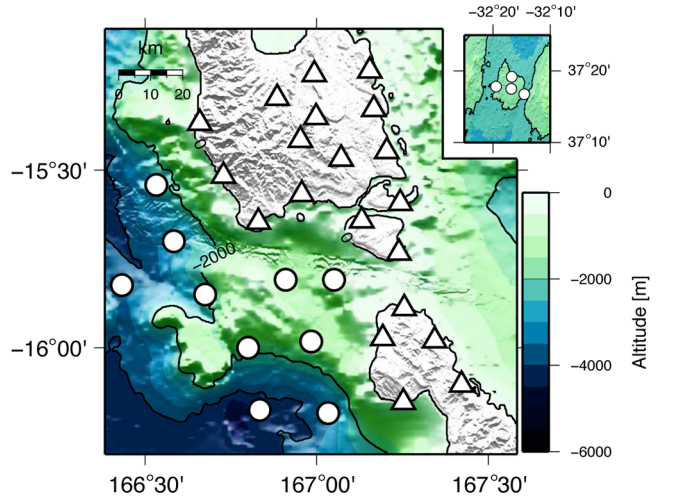
be long enough to correctly represent the onset energy and to avoid assigning high SNR to spikes. We estimate how well the onsets exceed the noise level by counting the number of times that the SNR function passes 20% (parameter  $T_{\text{SNR}}$ , Tables 1 and 2) of the maximum SNR. If this number exceeds a user-defined value  $v_c$  (usually 2, for  $P$  and  $S$  arrivals; Tables 1 and 2), the signal is rejected.

If the signal passes the above test, the algorithm picks up to two onsets using our CF. If the SNR value for either onset is below the user-defined threshold  $\text{SNR}_{\text{min}}$  (Tables 1 and 2), the corresponding pick is rejected. Finally, we apply the  $P$ – $S$  characterization procedure to the remaining onset(s) to classify them as  $P$  or  $S$  onsets or to reject them.

The APP then assigns a quality index to  $P$ - and  $S$ -wave onset picks based on their SNR. The indexes go from 0 (excellent quality pick) to 4 (rejected pick). These indexes are important because they define the weight of each pick in the inversion process used to calculate hypocenters. In order to give the  $P$  onsets a higher priority in the inversion, we set the SNR thresholds for  $S$  waves to 1.25 (parameter  $M_f$ , Tables 1 and 2) times the  $P$ -arrival SNR thresholds (parameters  $P_0$ ,  $P_1$ ,  $P_2$ , and  $P_3$ , Tables 1 and 2). This parameter depends mainly on the confidence we have in the  $P$  or  $S$  picking, but we find the value of 1.25 gives good results for land-based seismological datasets.

Once onsets are picked and characterized for each station, the APP applies a second rejection method, based on the distribution of onset times. This method is composed of two successive stages: elimination of strong outliers using clustering criteria, followed by removal of remaining outliers using statistical criteria. In the first stage, for each event, the APP sorts  $P$ -onset times in ascending order and creates clusters using the following definition: two consecutive onsets belong to the same cluster if and only if they are separated by less than a user-defined time lapse  $W_P$  (Tables 1 and 2). The APP rejects all onsets that do not belong to the biggest cluster (note if  $W_P$  is too big, only one cluster will be generated and no picks will be rejected). In the second stage, the APP calculates the median time of the remaining  $P$  onsets and the offset of each onset from this median, then rejects picks with offsets more than three standard deviations away from the median. It then applies the same method to  $S$  onsets (using a different time lapse [parameter  $W_S$ , Tables 1 and 2] in the clustering analysis). Finally, it applies only the second stage to  $P$ – $S$  delays (both  $P$  and  $S$  onsets are rejected in the case of  $P$ – $S$  delay outliers). This rejection method is appropriate for networks that are approximately uniformly distributed (the pick distributions follow a Gaussian law): a network with one site far away from the rest would have that site’s arrivals systematically rejected.

The final stage of the APP is to measure the amplitude and period of the peak arrival signal. The APP converts counts to displacement (using the seismometer plus digitizer transfer function) and applies a Wood–Anderson filter to all traces. The amplitude is the maximum from any channel for each station. The  $P$  and  $S$  offsets and the parameters  $A_0$  and



**Figure 7.** The seismic network geometries for Vanuatu (main map) and MAR (upper-right inset) datasets. The length and depth scales are the same for the two maps. The Vanuatu network is composed of 10 OBSs (circles) and 20 land stations (triangles). The MAR network is composed of four OBSs. The color version of this figure is available only in the electronic edition.

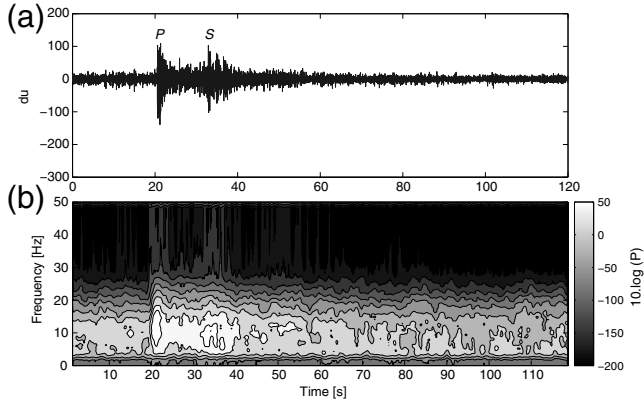
$A_{PS}$  (Tables 1 and 2) define the window in which to search for the maximum. If both  $P$  and  $S$  onsets are picked, the maximum amplitude is evaluated between the  $P$  pick and  $A_0$  seconds after the  $S$  pick; if only  $S$  is picked, the amplitude is evaluated between  $A_{PS}$  seconds before and  $A_0$  seconds after the  $S$  pick; and, if only  $P$  is picked, the amplitude is evaluated between the  $P$  pick and  $A_{PS} + A_0$  seconds after the  $P$  pick. The amplitude and its period are entered into the catalog so that the hypocenter inversion can estimate local magnitudes.

For each extracted event, the algorithm saves all information (station names, components, onset times, quality indexes, and maximum amplitudes and their periods) in a Nordic-format catalog file (Havskov and Ottemöller, 2010). This catalog file, plus a velocity model, can then be fed into one of many software codes to calculate event hypocenters and magnitudes.

### Application to the Vanuatu Dataset

We tested our method using data from a seismological network covering part of the Vanuatu subduction zone. The data were acquired in 2008–2009 as part of the Arc-Vanuatu experiment (see Data and Resources). The Vanuatu subduction zone is one of the most active seismic areas in the world, with more than 37 events of magnitude  $M_w \geq 7$  since 1973 (U.S. Geological Survey). The network was composed of 20 wideband seismometers onshore and 10 ocean-bottom seismometers. The data rate was 100 samples/s, the aperture was  $100 \times 100$  km, and the average distance between instruments was approximately 20 km; however, due to field constraints and the total absence of road access in the central part of the islands, the network is not as regular as one would





**Figure 8.** Frequency content of a typical seismogram from the Vanuatu dataset containing  $P$  and  $S$  onsets (station VANGO): (a)  $Z$  trace filtered between 3 and 45 Hz for clarity, and (b) spectrogram of the  $Z$  trace (grayscale gives power spectral density).

wish (Fig. 7). More than 100 events were detected per day; to cover a representative sample of local events, we chose all events from one day, and the only criterion we used is that the chosen day must have a good spatial distribution of events. We chose events from 1 June 2008, for which the event detector extracted 163 events.

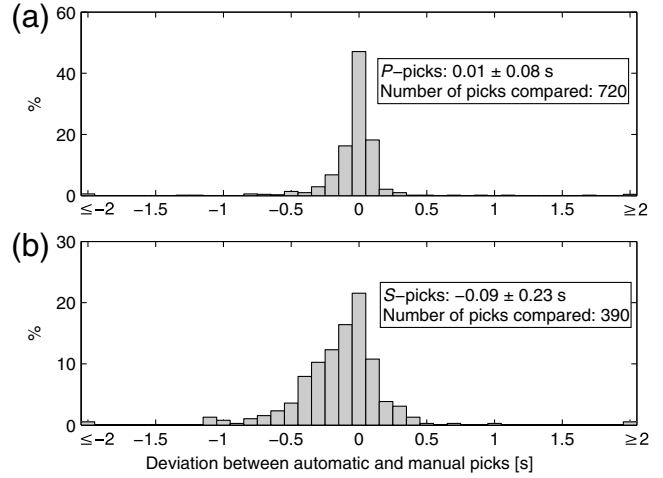
Before applying the APP, the frequency content of the seismogram should be evaluated, which will help to choose the set of frequency bandwidths used for the CF calculations. Figure 8 shows the spectrogram of a typical Vanuatu trace containing  $P$  and  $S$  onsets. The onset energy is concentrated between 3 and 18 Hz, which we used as our tightest frequency bounds (Table 2).

We were able to manually pick 99 of the 163 extracted events. Of these, 93 were locatable using the HYPOCENTER inversion program (Lienert *et al.*, 1986; Lienert and Havskov, 1995). Our APP returned 133 locatable events, including all 93 events located using the manual picks.

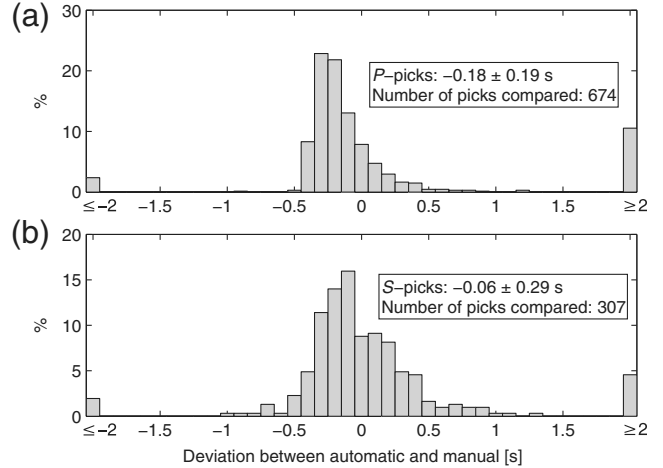
Manual picking provided 1094  $P$  and 507  $S$  onsets, whereas the APP provided 1007  $P$  and 1002  $S$  onsets. The automatic picker thus identified many more  $S$  onsets than the human operator, and about the same number of  $P$  onsets. However, the onsets picked were not always the same: 720 (66%) of the manual  $P$  onsets were picked automatically, as were 390 (77%) of the manual  $S$  onsets.

To evaluate the accuracy of the automatic picker, we analyze the time differences between the automatic and manual picks (Fig. 9). For  $P$  onsets, the median difference is  $0.01 \pm 0.08$  s (using the 68% interval as the variance; using the 95% interval, the variance would be  $\pm 0.40$  s). The small median and symmetrical variance distribution of time differences indicate there is little or no systematic shift with respect to manual picks. The kurtosis is known for its small systematic shifts compared with other metrics such as the skewness (Küperkoch *et al.*, 2010), which is one of the reasons we chose it for onset picking.

For  $S$  onsets, the median difference is  $-0.09 \pm 0.23$  s ( $\pm 0.61$  s for the 95% interval). The decrease in accuracy



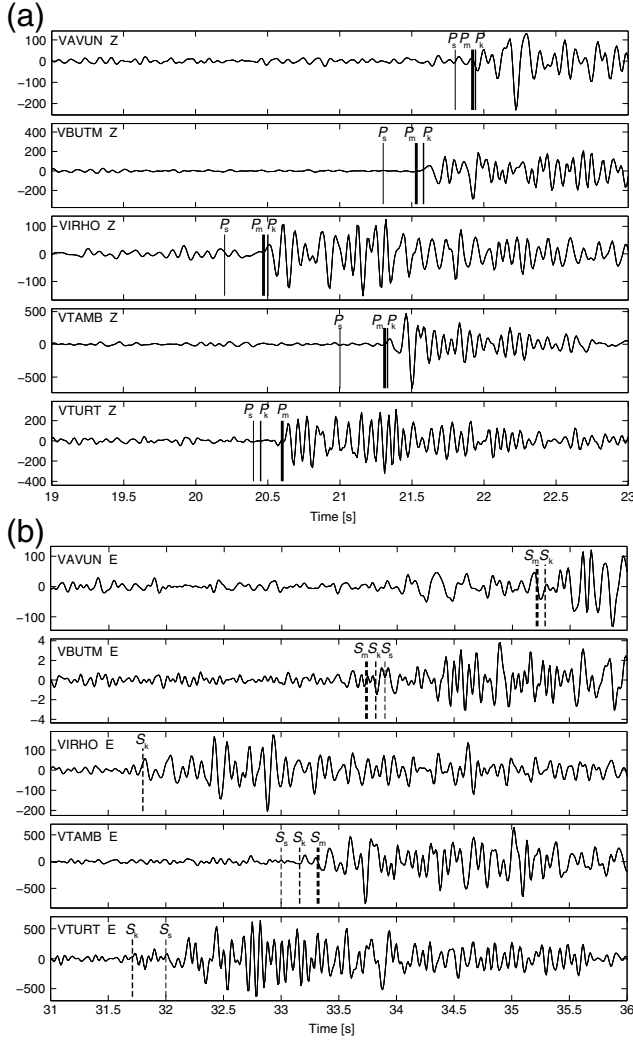
**Figure 9.** The difference between automatic and manual picking times using the proposed APP (Vanuatu dataset). All values beyond 2(−2) s are compiled at 2(−2) s. (a)  $P$ -wave onsets and (b)  $S$ -wave onsets.



**Figure 10.** The difference between automatic and manual picking times using the STA/LTA picking method (Vanuatu dataset). All values beyond 2(−2) s are compiled at 2(−2) s. (a)  $P$ -wave onsets and (b)  $S$ -wave onsets.

compared with  $P$ -wave picking is explained by the fact that  $S$  waves are generally emergent and buried in the  $P$  coda. The distribution shows a small trend to negative residuals, indicating that automatic picks are slightly ahead of manual picks. This may be because the kurtosis identifies the distribution change associated with an emergent onset before it is visible to the human eye.

We also applied an automatic STA/LTA picker to the data (LTA window = 10 s, STA window = 1 s). The STA/LTA picker made 1220  $P$  and 635  $S$  picks: 674 (62%) of the manual  $P$  onsets were picked, as were 307 (61%) of the manual  $S$  picks. The STA/LTA picker provided more  $P$  and less  $S$  onsets than our APP. For  $P$  onsets (Fig. 10a), the median difference between the STA/LTA and manual picks is  $-0.18 \pm 0.19$  s ( $\pm 15.60$  s for the 95% interval). The systematic offset is



**Figure 11.** Comparison between automatic picks (kurtosis-derived and STA/LTA pickers) and manual picks for  $P$  and  $S$  waves for one event of the Vanuatu project. Notations k, s, and m stand for kurtosis, STA/LTA, and manual picks. We selected five stations from our set of 30 stations. All traces shown are filtered between 3 and 45 Hz. (a)  $P$  onsets are shown on the  $Z$  component, and (b)  $S$  onsets are shown on the east component.

much larger than with our picking algorithm, and there are many more outliers, as indicated by the very large values for the 95% interval (in a Gaussian distribution, the time to the 95% interval would be twice that to the 68% interval). Outliers are mainly due to noise spikes that were picked by the STA/LTA picker. The comparison is not completely fair: our APP has stages dedicated to rejecting bad picks, whereas the STA/LTA picker does not. However, our APP is closer to manual picks even when we use metrics that ignore large mis-picks. For example, 46% of the compared  $P$  picks made by our APP are less than 0.1 s from the manual picks, compared with only 8% using the STA/LTA method. For  $S$  onsets (Fig. 10b), the median difference between the STA/LTA picks and manual picks is  $-0.06 \pm 0.29$  s ( $\pm 12.56$  s for the 95% interval): the

systematic offset is comparable to that for our automatic picker, but there are 38% less picks, and there are many more outliers. Using our APP, 22% of the common  $S$  picks are less than 0.1 seconds from the manual picks, compared with 9% using the STA/LTA method.

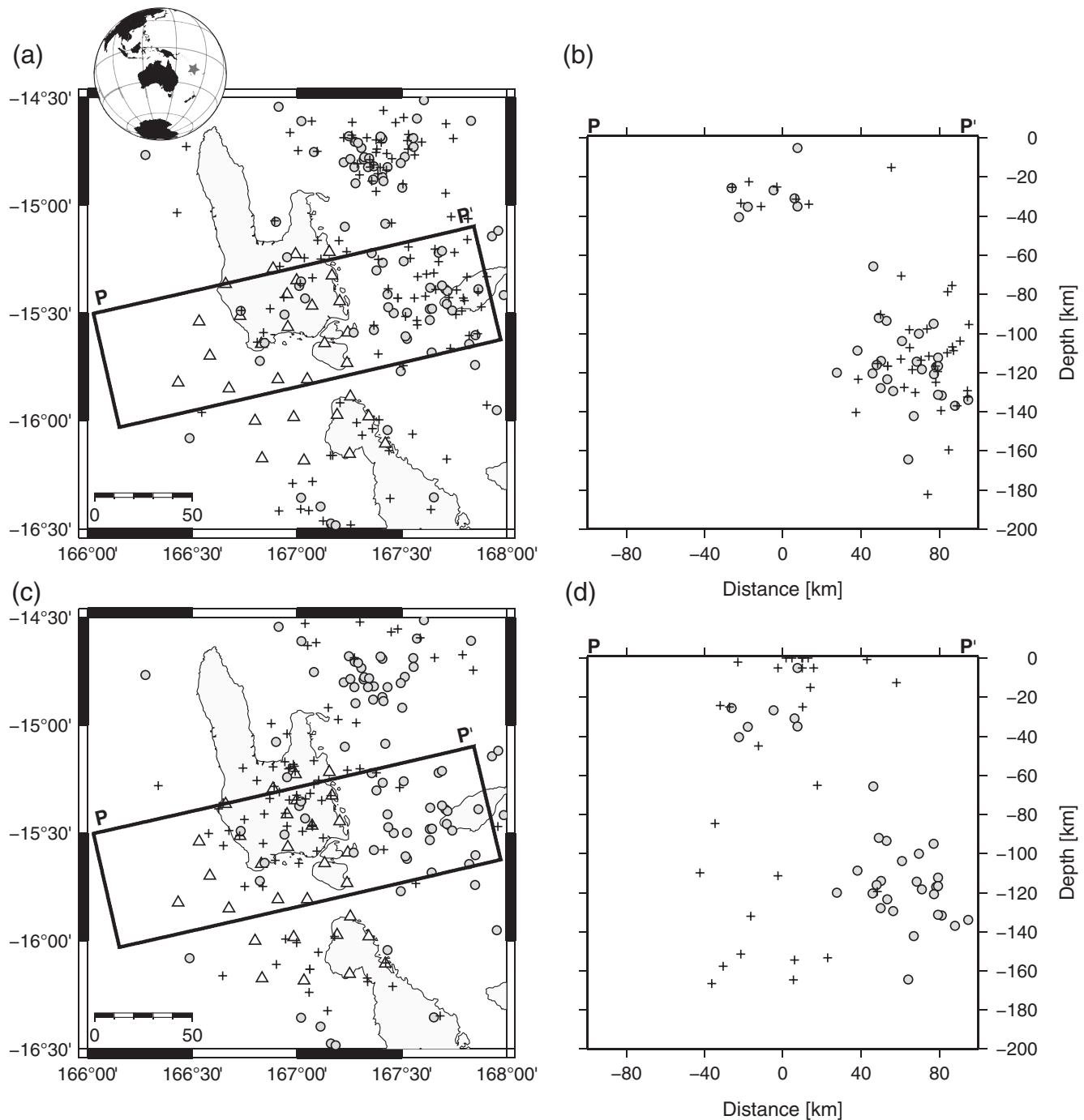
Figure 11 shows a typical example of manual and automatic picks for one event on traces from five network stations. The time differences between automatic and manual picks are smaller using our picker than using the STA/LTA picker, both for high SNR ( $P$  onsets, Fig. 11a) and small SNR ( $S$  onsets, Fig. 11b).

Figure 12 compares hypocenters calculated using onset times provided by the three different picking methods. The APP hypocenters are much closer to the manual hypocenters (Fig. 12a,b) than are the STA/LTA hypocenters (Fig. 12c,d).

We now evaluate the reliability of the quality indexes assigned by our picking method. We compare automatically and manually assigned quality indexes, then look at how the time differences between automatic and manual picks change as a function of the quality index. Figure 13 shows a histogram of manual quality indexes subtracted from automatic quality indexes, for both  $P$  and  $S$  onsets. If the analyst and the algorithm assigned the same quality index to every pick, the histogram would be a single bar at 0. Instead, the histograms are distributed on an approximately Gaussian shape and their centers are offset from zero: the automatically assigned  $P$ -onset quality indexes are slightly negative and the  $S$ -onset quality indexes are slightly positive when compared with the manual values. The offset is deliberate and comes from our choice of SNR thresholds and  $M_f$  (Table 2); we want the hypocenter location routine to give more weight to  $P$ - than to  $S$ -wave onsets, and we had to slightly improve the  $P$  qualities to avoid rejecting  $S$  onsets. The user can assign his or her own values to obtain the desired center points by changing the SNR values and  $M_f$ . The distribution around the central value is not perfect (assuming all of the manual quality indexes were perfectly assigned), but 85% of the  $P$  picks and 75% of the  $S$  picks are within 1 of the central value.

We now look at the relationship between the picking accuracy of our algorithm and the automatically assigned quality indexes. For both  $P$  and  $S$  residuals, the standard deviation tends to increase for bigger quality indexes (Fig. 14), which is consistent with these bigger numbers corresponding to lower pick certainty. The differences are not as large as we might expect if we consider that an increase of 1 in the weight number is generally treated as a factor of 2 decrease in the pick time certainty, but we are comparing between picks and not against true arrival times. The relatively small differences between the residuals as a function of quality index probably indicate the APP onset estimates are consistent with the manual estimates. Table 3 summarizes the distribution of residuals and their variances as a function of the quality indexes.

We also evaluate the relationship between picking accuracy and earthquake local magnitudes (calculated from our

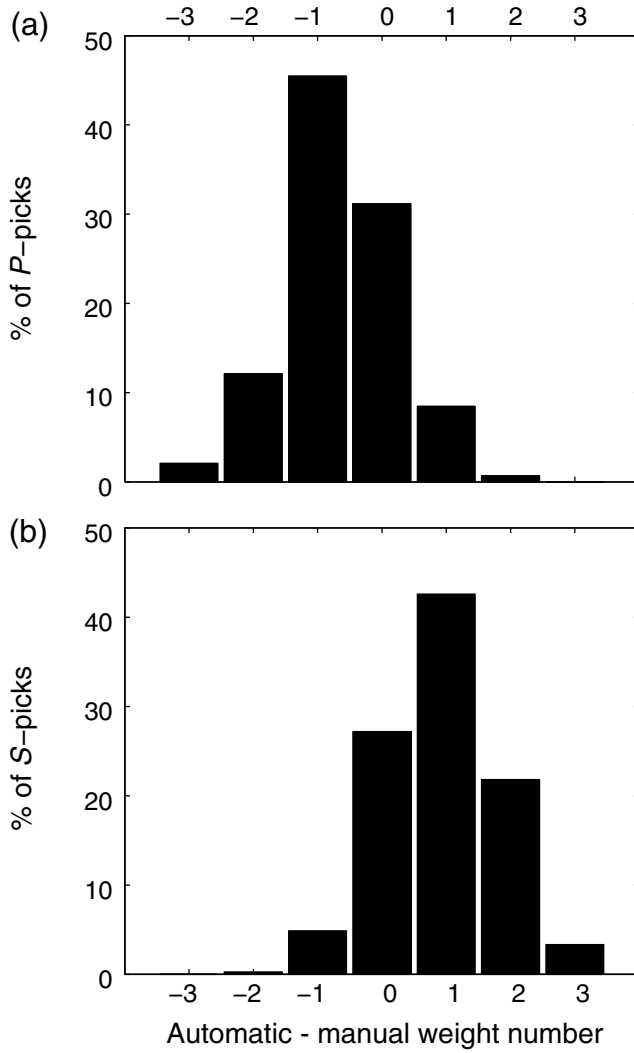


**Figure 12.** Map and profile views showing hypocenters resulting from automatic (kurtosis-derived and STA/LTA) and manual picking for the Vanuatu region. Triangles, stations; circles, hypocenters derived from manual picking; and crosses, hypocenters derived from automatic picking. (a,b) Map and profile view (using  $P$ - $P'$  projection) comparing hypocenters localized using manual and kurtosis-derived picks. (c,d) Same as (a,b), comparing manual and STA/LTA-derived hypocenters.

automatic amplitude estimates). Figure 15a shows the distribution of local magnitudes of the events we used. The events magnitudes are between 0.7 and 3.7, with half of the events below magnitude 1.9. There is no clear correlation between magnitudes and the residuals or their variances (Fig. 15b-e), indicating that the automatic picker's accuracy (relative to manual picking) is independent of the event magnitude.

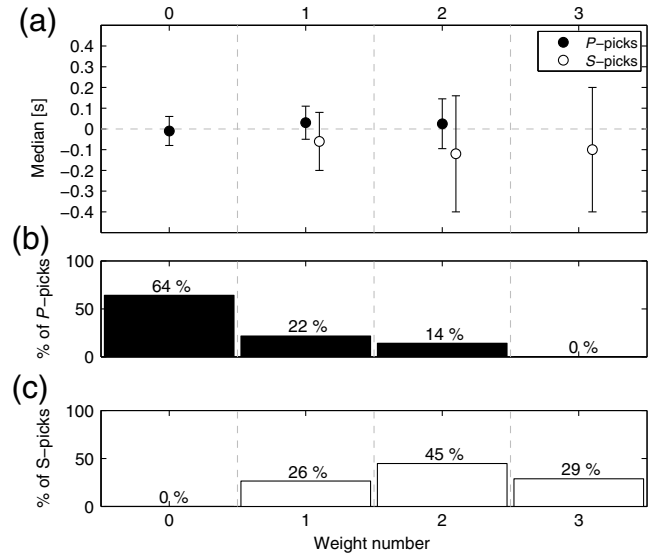
#### Application to the Mid-Atlantic Ridge Dataset

We next tested our method on a very different dataset: four ocean-bottom seismometers around a seafloor volcano (Fig. 7 inset). The network had an aperture of 7 km, the sensors had a fourth channel ( $H$ ) for pressure data (measured by



**Figure 13.** Histograms of the difference between pick quality indexes assigned automatically and manually (automatic minus manual) for (a)  $P$  onsets, and (b)  $S$  onsets.

a hydrophone), and the seismometer data were short period (4.5 Hz corner frequency). The events within the network had magnitudes from  $-1.2$  to  $1.5$  and depths from 2 to 3.5 km beneath the seafloor (Crawford *et al.*, 2013). As op-



**Figure 14.** Relationship between the automatically assigned quality indexes and the pick quality. (a) Differences between automatic and manual picks and their standard deviations. (b) The distribution of quality indexes for  $P$ -onset picks. (c) The distribution of quality indexes for  $S$ -onset picks.

posed to the Vanuatu data, MAR  $P$  onsets do not have a clear signature in the frequency domain (Fig. 16).

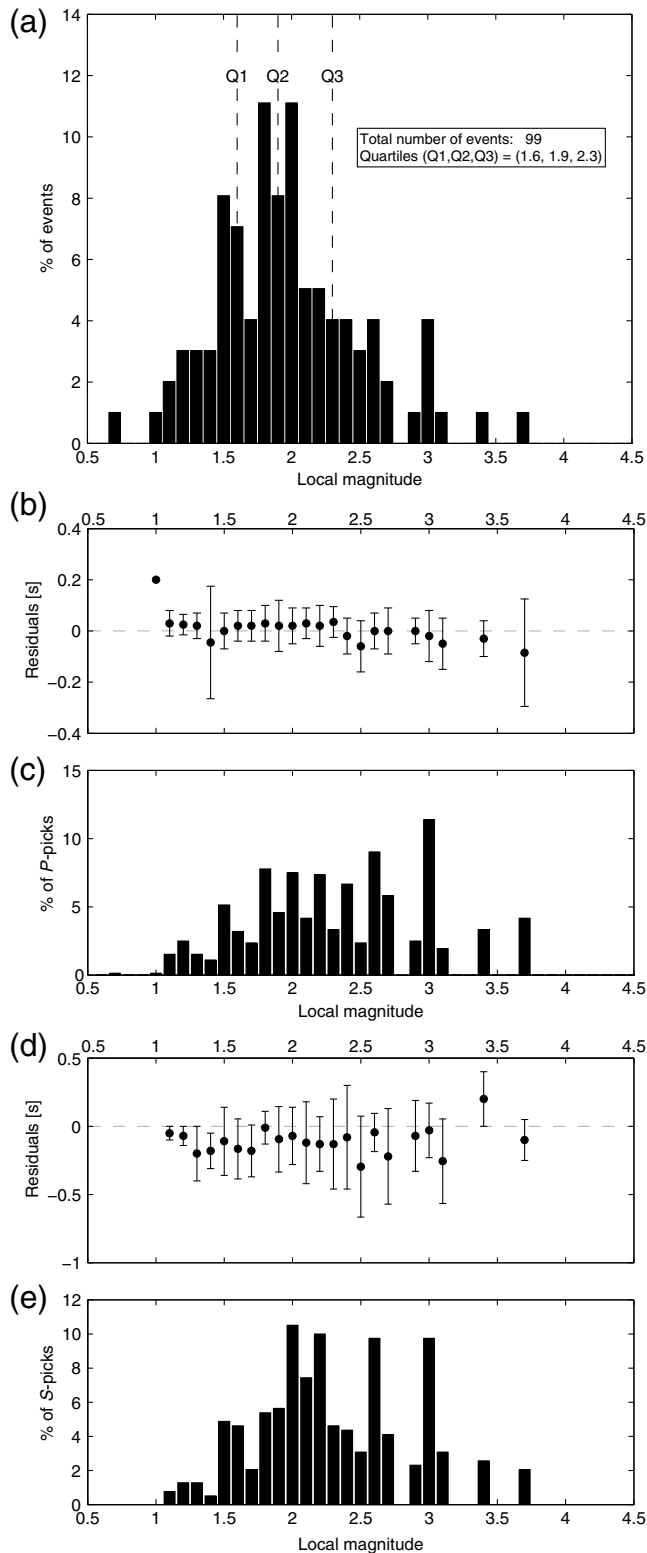
We automatically picked one month of data (February 2011) that had been manually picked. These data are particularly difficult to analyze automatically for several reasons: (1)  $P$  onsets were generally much less energetic than  $S$  onsets and were clear only on the  $H$  component; (2) sea-surface reflected  $P$  waves arrived between the  $P$  and  $S$  onsets, complicating pick selection and identification; (3) some components were unusable, so we could not apply the polarity approach; (4) the small number of instruments required us to be very careful in our outlier rejection in order to retain as many valid picks as possible.

To respond to these difficulties, we applied some changes to the APP. (1) The APP only computed the CF for the best component of each station. (2) The  $S$  onset is identified by taking the onset located just before the maximum trace energy (energy conditioning). (3) It then identified the

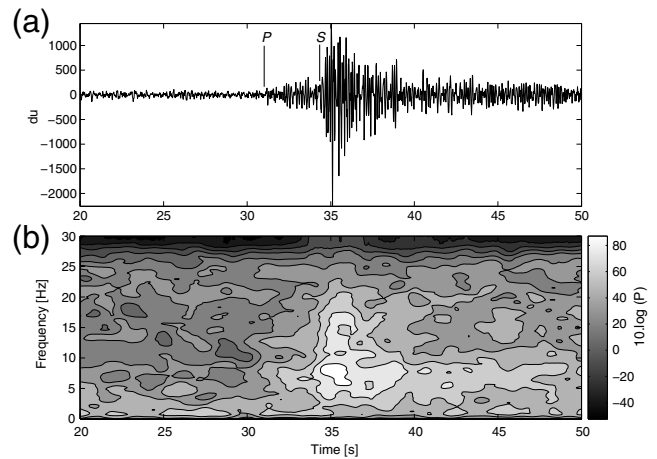
Table 3  
Differences between Manual  $P$  and  $S$  Picks for Both APP and STA/LTA Picks for the Vanuatu Dataset

Method	Weight	$P$ Picks				$S$ Picks			
		Fraction (%)	Residuals (s)	Variance		Fraction (%)	Residuals (s)	Variance	
				$\sigma$	$2\sigma$			$\sigma$	$2\sigma$
New Kurtosis	0	64	-0.01	0.07	-	0	-	-	-
	1	22	0.03	0.07	-	26	-0.06	0.16	-
	2	14	0.03	0.10	-	45	-0.12	0.23	-
	3	0	-	-	-	29	-0.10	0.28	-
	All	100	0.01	0.08	0.40	100	-0.09	0.23	0.61
STA/LTA	All	100	-0.18	0.19	15.60	100	-0.06	0.29	12.56

For the APP picks, these differences and the percentage of picks assigned are also shown as a function of the quality index.



**Figure 15.** Relation between picking accuracy and local magnitude. (a) Magnitude distribution of 99 compared events (picked manually and picked and located by the automatic picker). The bin size is 0.1. Quartiles 1, 2, and 3 indicate the values below which are 25%, 50%, and 75% of the events, respectively. (b) Residuals and standard deviations for  $P$ -onset picks as a function of magnitude. (c) The distribution of  $P$  picks as a function of magnitude. (d,e) Same as (b) and (c), but for  $S$ -onset picks.



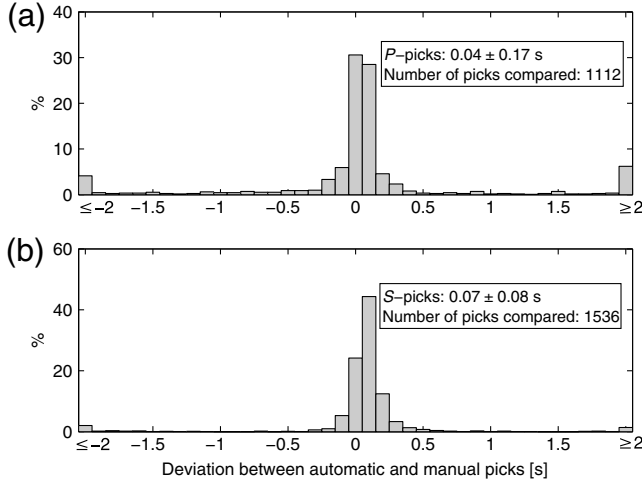
**Figure 16.** Frequency content of a typical seismogram from the MAR dataset containing  $P$  and  $S$  onsets (station LSd1). (a)  $Z$  trace filtered between 5 and 18 Hz for clarity. (b) Trace's spectrogram in which intensities refer to power spectral density.

$P$  onset using an algorithm explained in the next paragraph. Choice of components to be processed and energetic conditioning are user-defined options that can be directly defined in the algorithm input file. All user-defined parameters used in the APP for the MAR dataset are summarized in Table 2.

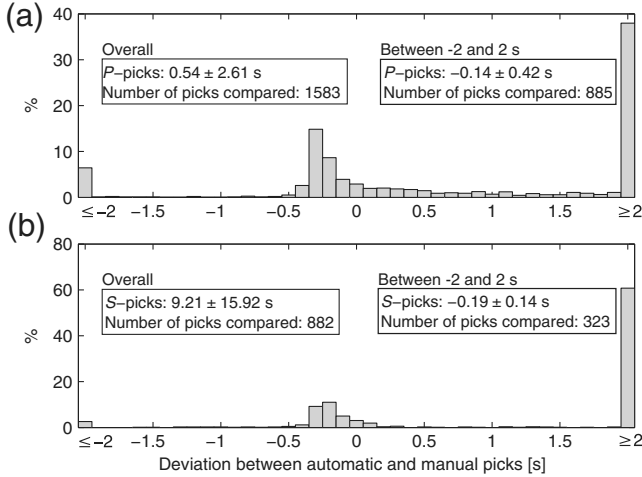
In the marine data, to identify the  $P$  onset, the APP first selects the strongest minimum of the CF preceding the  $S$  pick. This is the  $P$ -onset candidate. It then checks if there is another strong minimum at a time  $\Delta t^P = (2H)/V_w$  before the  $P$  onset candidate, in which  $H$  is the water depth in meters and  $V_w$  is the estimated acoustic-wave velocity in water ( $\approx 1500$  m/s). If there is, the  $P$ -onset candidate was actually a sea-surface bounce, and the APP selects the preceding strong minimum as the true  $P$  onset.

Manual picking provided 1801  $P$  and 1809  $S$  onsets, whereas the APP provided 1676  $P$  and 2349  $S$  onsets. We compared automatic and manual picks (Fig. 17). 1112 (61%) of the manual  $P$  onsets and 1536 (84%) of the manual  $S$  onsets were picked automatically. For  $P$  waves, the median difference is  $0.04 \pm 0.17$  s, for  $S$  waves it is  $0.07 \pm 0.08$  s. A smaller standard deviation is obtained for  $S$  than for  $P$  picks because of the high energy of  $S$  onsets.

We also automatically picked the data using the STA/LTA picker (Fig. 18). In this case median differences and standard deviations are severely biased by the large amount of outliers (37% of the  $P$  picks residuals and 69% of the  $S$  residuals are above 2 s). These outliers are in part caused by the STA/LTA picker picking  $S$  onsets as  $P$  onsets and spikes as  $S$  onsets. If we do not consider outliers (residuals above 2 s in absolute value), only 49% of the manual  $P$  and 17% of the manual  $S$  onsets were picked by the STA/LTA. Compared to our picker, the STA/LTA causes much more mis-picks. When throwing out the outliers the median and the residual are  $-0.14 \pm 0.42$  s for  $P$  onsets and  $-0.19 \pm 0.14$  s for  $S$  onsets. Residuals obtained by the two methods are summarized in Table 4. Figure 19 shows manual and automatic

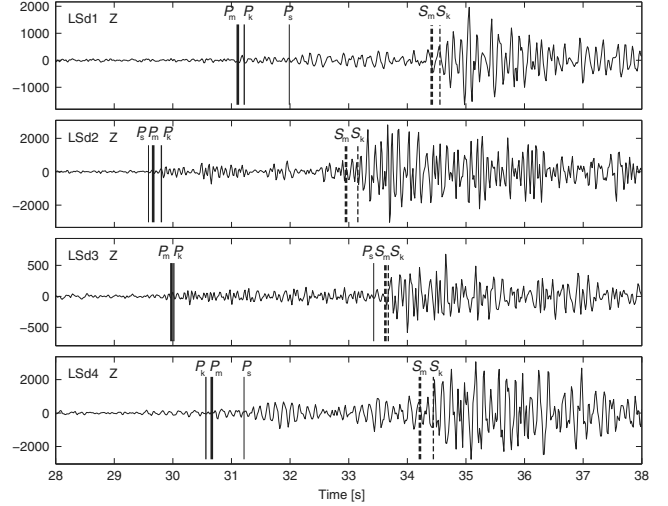


**Figure 17.** The difference between automatic and manual picking times using the proposed APP (MAR dataset). All values beyond 2(−2) seconds are compiled at 2(−2) seconds for (a)  $P$ -wave onsets and (b)  $S$ -wave onsets.



**Figure 18.** The difference between automatic and manual picking times using the STA/LTA picking method (MAR dataset). All values beyond 2(−2) seconds are compiled at 2(−2) seconds. Median and standard deviations are shown for overall picks and for residuals in between −2 and 2 s. (a)  $P$ -wave onsets. (b)  $S$ -wave onsets.

picks on four vertical channel traces of our network for one event. Again, picks from the kurtosis-derived picker are closer to manual picks than the ones from the STA/LTA picker. Moreover in this and many other cases, the STA/LTA



**Figure 19.** Comparison between automatic (kurtosis-derived and STA/LTA) and manual picks for  $P$  and  $S$  onsets for one event of the MAR dataset. Indexes under  $P$  and  $S$  are the same as in Figure 11. We plotted the  $Z$  component of all four OBSs of the network. All traces shown are filtered between 3 and 20 Hz.

picker has difficulty picking  $S$  waves and even declares an  $S$  onset as a  $P$  onset (OBS station LSD3).

These results demonstrate that the APP with our new CF can be applied to a wide range of networks and instrument types. The APP has also been adapted to localize rock fall events occurring in the crater of the Piton de la Fournaise volcano on Réunion island (Hibert *et al.*, 2011; Hibert, 2012). Picking shows good performance even when dealing with emergent waves recorded on a few number of stations (seismic network composed of five stations).

## Discussion

In the past two decades, many new automatic pickers have been proposed, using different approaches and methods. Even if significant improvements have been made, only a few of these algorithms are openly available to the scientific community and can be rapidly implemented on new datasets. Two of these are included in the well-known Earthworm (Johnson *et al.*, 1995) and SeisComp3 (Hanka *et al.*, 2000) software suites. The automatic picking module implemented in both software suites is based on the STA/LTA method introduced by Allen (1982), which is simple, robust, and easily adaptable. Implementation of kurtosis-derived

Table 4  
Differences between Manual  $P$  and  $S$  Picks and Both APP and STA/LTA Picks, for the MAR Dataset

Method	Comments	$P$ Picks		$S$ Picks	
		Residuals (s)	Variance $\sigma$	Residuals (s)	Variance $\sigma$
New Kurtosis	All	0.04	0.17	0.07	0.08
STA/LTA	All	0.54	2.61	9.21	15.92
	$-2 \leq \text{residuals (s)} \leq 2$	-0.14	0.42	-0.19	0.14

methods in these suites should allow much more accurate onset picking.

Two other well-known picking algorithms are the “tandem picker” (Nippres *et al.*, 2010) and the MannekenPix (MPX) algorithm (Aldersons, 2004; Di Stefano *et al.*, 2006). Both methods require operator-intensive preparation phases, and neither automatically identifies phase arrivals. The MPX method uses the Baer–Kradolfer algorithm (Baer and Kradolfer, 1987), which has been shown to be less accurate than the kurtosis (Küperkoch *et al.*, 2010), whereas the tandem picker uses the kurtosis but picks at the point of maximum slope, which we have shown to be less accurate than picking the initial point of inflection.

AR-AIC methods are powerful tools to pick onset times when the SNR is very low (Takanami and Kitagawa, 1993; Leonard and Kennett, 1999), but their implementation is very computer intensive, as a large number of AR models must be calculated. Taking into account these considerations, we believe the kurtosis-derived method that we developed here is the most appropriate tool for fast and automatic onset picking.

The dip–rectilinearity polarization analysis proposed in this paper is well suited for independently identifying *P*, *S*, and surface waves, in most cases. However, the polarization analysis does not reliably identify *P* and *S* arrivals if the *P*–*S* delay is shorter than the analysis window. We demonstrated that this window should be longer than three cycles of the lowest period of the arrival waveform: if there is a risk of significant arrivals with shorter *P*–*S* delays, it is recommended that the user run the algorithm twice—once using the polarization analysis and once without—to see if significant events are lost using the polarization analysis. If this is the case, the extra events identified without the polarization analysis could be manually added into the database. The non-polarization analysis code identifies the phase of events using energetic conditioning (used in the MAR example), which is less discriminating than polarization analysis but which does not rely on an analysis window length.

A future prospect could be to automatically switch from polarization analysis to another method (such as energetic conditioning) when the time delay between identified onsets is shorter than the polarization window length, or to output undefined arrivals with an appropriate flag.

## Conclusion

We have presented a new automated algorithm to identify and pick *P*- and *S*-wave onsets. The method uses a new CF based on the kurtosis statistical function to accurately pick seismic onsets. The CF uses several sliding window lengths and multiple frequency bandwidths, minimizing the dependence of the picker on these parameters. A polarization analysis is applied to distinguish *P* waves from *S* waves and to reject picks of surface waves or noise. A second rejection method is then applied, based on the clustering of *P*-onset times, *S*-onset times, and *P*–*S* onset time differences. The algorithm includes a pick quality classification based on

the SNR energy ratio, and it calculates amplitudes to allow magnitude estimates.

We tested the automatic picker’s performance using two datasets: a 30-instrument land-based array with wideband seismometers, and a four-instrument seafloor array—often without three-component data—that pushes the limits for automatic picking and event location. For the first network, of 163 automatically selected events, 99 events could be picked manually, and 93 of these could be located, whereas the automatic picker provided 133 locations. The automatic picker picked as many *P* onsets and more *S* onsets than manual picking, overlapping the manual onsets on 66% of the *P* picks and 77% of the *S* picks. The automatic picks deviated from the manual picks by  $0.01 \pm 0.08$  s for *P* and  $-0.09 \pm 0.23$  s for *S*. The offset and deviation are much better than obtained using an STA/LTA picker and, in addition, many more *S* arrivals were picked. The assigned pick quality index, based on the SNR, correlates well with manual quality indexes. The pick accuracy (relative to manual picks) is independent of the event magnitude.

The second dataset confirms the improved performance of our APP compared with the STA/LTA method, for which over 40% of both *P* and *S* picks were more than 2 s from the manual picks.

The automatic picking algorithm described in this paper can be a powerful tool for automatically picking *P* and *S* onsets with high precision and accuracy and coherently assigning their quality index. Very few manually picked events were lost and, for the tested datasets, the good quality and consistency of the picking allowed more events to be located. The number of picked events and accuracy of the picking are significantly higher using our APP than using the STA/LTA method. The picker has few and relatively simple user-defined parameters and should be easily adaptable to a wide range of local networks.

## Data and Resources

Data used in this paper come from projects funded by IRD through Géoazur and by the French Ministry of Research through the ANR Arc-Vanuatu program (Vanuatu data) and by MoMAR and EMSO Azores (MAR data). Data were analyzed using MATLAB 2012b (<http://www.mathworks.fr/products/matlab/>) and SEISAN (<http://seis.geus.net/software/seisan/>, last accessed on March 2013).

## Acknowledgments

We would like to thank Alexandre Necessian, Francois Beauducel, Claudio Satriano, and Pascal Bernard from the Institut de Physique du Globe de Paris (IPGP) seismological laboratory for their advice and contributions to this paper. The Vanuatu data were collected as part of the ARC-Vanuatu program funded by the French National Research Agency (ANR). We would like to thank Marc Régnier for his contribution to this project. The MAR data were collected as part of the French MoMAR and EMSO Azores initiatives, funded by l’Institut Français de Recherche pour l’Exploitation de la Mer (IFREMER), l’Institut National des Sciences de l’Univers (INSU), and the Centre National de Recherche Scientifique (CNRS).

We also thank Gilles Grandjean from the Bureau de Recherches Géologiques et Minières (BRGM) and the ANR UNDERVOLC program. This is IPGP contribution 3420.

## References

- Akaike, H. (1974). Markovian representation of stochastic processes and its application to the analysis of autoregressive moving average processes, *Ann. Inst. Stat. Math.* **26**, 363–387.
- Aldersons, F. (2004). Toward a three-dimensional crustal structure of the Dead Sea region from local earthquake tomography, *Ph.D. Thesis*, Tel Aviv University, Israel, 120 pp., <http://faldersons.net> (last accessed June 2013).
- Allen, R. (1982). Automatic phase pickers: Their present use and future prospects, *Bull. Seismol. Soc. Am.* **72**, S225–S242.
- Baer, M., and U. Kradolfer (1987). An automatic phase picker for local and teleseismic events, *Bull. Seismol. Soc. Am.* **77**, 1437–1445.
- Colaço, A., J. Blandin, M. Cannat, T. Carval, V. Chavagnac, D. Connelly, M. Fabian, S. Ghiron, J. Goslin, and J. M. Miranda (2011). MoMAR-D: A technological challenge to monitor the dynamics of the Lucky Strike vent ecosystem, *ICES J. Mar. Sci. (J. du Conseil)* **68**, 416–424.
- Crawford, W. C., R. Abhishek, S. C. Singh, M. Cannat, J. Escartin, H. Wang, R. Daniel, and V. Comber (2013). Hydrothermal seismicity beneath the summit of Lucky Strike volcano, Mid-Atlantic Ridge, *Earth Planet Sci. Lett.* **373**, 118–128.
- DeCarlo, L. T. (1997). On the meaning and use of kurtosis, *Psychol. Meth.* **2**, 292–307.
- Di Stefano, R., F. Aldersons, E. Kissling, P. Baccheschi, C. Chiarabba, and D. Giardini (2006). Automatic seismic phase picking and consistent observation error assessment: Application to the Italian seismicity, *Geophys. J. Int.* **165**, 121–134.
- Freedman, H. W. (1966). The “little variable factor” a statistical discussion of the reading of seismograms, *Bull. Seismol. Soc. Am.* **56**, 593–604.
- Gomberg, J. S., K. M. Shedlock, and S. W. Roecker (1990). The effect of *S*-wave arrival times on the accuracy of hypocenter estimation, *Bull. Seismol. Soc. Am.* **80**, 1605–1628.
- Hanka, W., A. Heinloo, and K. H. Jaeckel (2000). Networked seismographs: GEOFON real-time data distribution, *ORFEUS Electronic Newsl.* **2**.
- Havskov, J., and L. Ottemöller (2010). *Routine Data Processing in Earthquake Seismology: With Sample Data, Exercises and Software*, Springer, New York, 326 pp.
- Hibert, C. (2012). L'apport de l'écoute sismique pour l'étude des éboulements du cratère Dolomieu, Piton de la Fournaise, La Réunion, *Ph.D. Thesis*, Institut de Physique du Globe de Paris, France (in French).
- Hibert, C., G. Grandjean, N. Shapiro, and C. Baillard (2011). Automatic picking and localization of rockfalls from their seismic signature: Application to the Piton de la Fournaise volcano, La Réunion, poster, *AGU Fall Meeting 2011*.
- Hildyard, M. W., S. E. J. Nippress, and A. Rietbrock (2008). Event detection and phase picking using a time-domain estimate of predominate period Tpd, *Bull. Seismol. Soc. Am.* **98**, 3025–3032.
- Johnson, C. E., A. Bittenbinder, B. Bogaert, L. Dietz, and W. Kohler (1995). Earthworm: A flexible approach to seismic network processing, *IRIS Newsl.* **14**, 1–4.
- Jurkevics, A. (1988). Polarization analysis of three-component array data, *Bull. Seismol. Soc. Am.* **78**, 1725–1743.
- Küperkoch, L., T. Meier, J. Lee, and W. Friederich (2010). Automated determination of *P*-phase arrival times at regional and local distances using higher order statistics, *Geophys. J. Int.* **181**, 1159–1170.
- Leonard, M., and B. L. N. Kennett (1999). Multi-component autoregressive techniques for the analysis of seismograms, *Phys. Earth. Planet. In.* **113**, 247–263.
- Lienert, B. R., and J. Havskov (1995). A computer program for locating earthquakes both locally and globally, *Seismol. Res. Lett.* **66**, 26–36.
- Lienert, B. R., E. Berg, and L. N. Frazer (1986). Hypocenter: An earthquake location method using centered, scaled, and adaptively damped least squares, *Bull. Seismol. Soc. Am.* **76**, 771–783.
- Nippress, S. E. J., A. Rietbrock, and A. E. Heath (2010). Optimized automatic pickers: Application to the ANCORP data set, *Geophys. J. Int.* **181**, 911–925.
- Pelletier, B., S. Calmant, and R. Pillet (1998). Current tectonics of the Tonga–New Hebrides region, *Earth Planet. Sci. Lett.* **164**, 263–276.
- Saragiotis, C. D., L. J. Hadjileontiadis, and S. M. Panas (2002). PAI-S/K: A robust automatic seismic *P* phase arrival identification scheme, *IEEE Trans. Geosci. Remote Sens.* **40**, 1395–1404.
- Sleeman, R., and T. van Eck (1999). Robust automatic *P*-phase picking: An on-line implementation in the analysis of broadband seismogram recordings, *Phys. Earth. Planet. In.* **113**, 265–275.
- Takanami, T., and G. Kitagawa (1993). Multivariate time-series model to estimate the arrival times of *S*-waves, *Comput. Geosci.* **19**, 295–301.
- Vidale, J. E. (1986). Complex polarization analysis of particle motion, *Bull. Seismol. Soc. Am.* **76**, 1393–1405.
- Zeiler, C., and A. A. Velasco (2009). Seismogram picking error from analyst review (SPEAR): Single-analyst and institution analysis, *Bull. Seismol. Soc. Am.* **99**, 2759–2770.

Université Paris Diderot–PRES Sorbonne Paris Cité  
 Institut de Physique du Globe de Paris  
 UMR CNRS 7154  
 1 rue Jussieu  
 75238 Paris, France  
 baillard@ipgp.fr  
 crawford@ipgp.fr  
 hibert@ipgp.fr  
 mangeney@ipgp.fr  
 (C.B., W.C.C., C.H., A.M.)

Université de La Rochelle  
 Littoral, Environnement et Sociétés (LIENSIS)  
 UMR CNRS 7266  
 2 rue Olympe de Gouges  
 17000 La Rochelle, France  
 valerie.ballu@univ-lr.fr  
 (V.B.)

Manuscript received 29 November 2012