

A frame selective dynamic programming approach for noise robust pitch estimation

Chiranjeevi Yarra,^{1, a)} Om D. Deshmukh,² and Prasanta Kumar Ghosh¹

¹*Department of Electrical Engineering, Indian Institute of Science (IISc),
Bangalore, 560012, India*

²*Xerox Research Center India, Bangalore, 560103, India*

1 The principles of the existing pitch estimation techniques are often different and com-
2 plementary in nature. In this work, we combine the complementary characteristics of
3 two existing methods, namely, sub-harmonic to harmonic ratio (SHR) and sawtooth-
4 wave inspired pitch estimator (SWIPE), to improve pitch estimation. Using variants
5 of SHR and SWIPE, the proposed method, named FSDP, classifies all the voiced
6 frames into two classes – the first class consists of the frames where a confidence
7 score maximization criterion is used for pitch estimation, while for the second class,
8 a dynamic programming (DP) based approach is proposed. Experiments are per-
9 formed on speech signals separately from KEELE, CSLU and PaulBaghsaw corpora
10 under clean and additive white Gaussian noise at 20, 10, 5, and 0dB SNR conditions
11 using four baseline schemes including SHR, SWIPE and two DP based techniques.
12 The pitch estimation performance of FSDP, when averaged over all SNRs, is found
13 to be better than those of the baseline schemes suggesting the benefit of applying
14 smoothness constraint using DP in selected frames in the proposed FSDP scheme.
15 The VuV classification error from FSDP is also found to be lower than that from all
16 four baseline schemes in almost all SNR conditions on three corpora.

^{a)} chiranjeeviy@iisc.ac.in

17 I. INTRODUCTION

18 Accurate estimation of pitch is useful in various applications including gender classification¹,
19 emotion recognition², automatic intonation identification³, automatic music transcription⁴,
20 query by humming⁵, speech disorders identification⁶ and source-filter model based speech
21 coding systems^{7,8}. The reliability of these applications depends on the accuracy of the pitch
22 estimation. Typically pitch is considered as the fundamental frequency of the quasi-periodic
23 speech signal perceived by the human auditory system⁹⁻¹³. An accurate pitch estimation for
24 speech signal is non-trivial because – a) speech is not perfectly periodic due to non-stationary
25 variations in the frequency and the amplitude⁸, b) speech can be noisy, for example, in the
26 case where the distance between the microphone and speaker is large, c) the signal to noise
27 ratio (SNR) can be low¹⁴.

28 In the literature, several estimation techniques¹⁵⁻¹⁷ have used the dynamic programming
29 (DP) to impose temporal continuity in the estimated pitch contour^{12,18-24}. Typically, DP
30 based approaches divide the input signal into frames and identify multiple pitch candidates
31 for every frame. Often, these candidates are associated with a measure of confidence¹⁸,
32 referred to as confidence-score. These scores are considered in DP cost function for selecting
33 the best candidate in a frame, which is declared as the estimated pitch at that frame.

34 The manner in which the pitch candidates and their confidence-scores are computed
35 varies across different DP approaches. For example, in neural network based approaches,
36 probabilistic outputs of pitch candidate states are produced where a typical number of pitch
37 candidates is approximately 68^{22,23}. Among these, a deep neural network (DNN) based

38 approach is shown to be effective in noisy speech, although it requires a large amount of
39 data for training. In contrast to these data driven approaches, several knowledge based
40 approaches are proposed with fewer number of pitch candidates. The robust algorithm for
41 pitch tracking (RAPT) uses normalized cross correlation function (NCCF) to estimate mul-
42 tiple pitch candidates and their confidence-scores²⁰. The yet another algorithm for pitch
43 tracking (YAAPT) also uses NCCF to estimate multiple pitch candidates, but these candi-
44 dates are further refined using spectral information²⁵. The algorithm proposed by Ba et al.,
45 named BaNa, computes multiple pitch candidates by combining the approaches of harmonic
46 ratio and cepstrum analysis²⁶. The algorithm proposed by Gonzalez et al., named PEFA,
47 uses a convoluted normalized periodogram to estimate multiple candidates and then pitch
48 is estimated using DP¹². While RAPT and YAAPT have been shown to perform well for
49 clean and telephone channel speech respectively, PEFAC has been shown to perform better
50 in low SNR conditions. However at higher SNRs and clean conditions, PEFAC does not
51 have a satisfactory performance. In this work we propose a frame selective DP (FSDP)
52 approach which works better with few number of pitch candidates in clean as well as in
53 noisy conditions in both high and low SNR conditions. The proposed FSDP exploits speech
54 characteristics in order to estimate pitch using small amount of training data.

55 Similar to the pitch candidates, the computation of the confidence-scores plays an impor-
56 tant role in pitch estimation performance for both clean and noisy conditions. For example,
57 RAPT has a robust confidence-score computation associated with each candidate, which
58 could be the reason for it to have a better accuracy in clean case compared to other DP
59 based algorithms. However RAPT involves the selection of many parameters on a training

60 corpus, causing performance degradation across corpora as well as in noisy conditions¹².
61 To improve the performance under noisy conditions, PEFAC introduces another confidence-
62 score computation that uses several parameters heuristically designed during training under
63 noisy conditions.

64 Another critical factor for the performance of a DP based method is the weight given
65 to the continuity constraint. While the continuity constraint often helps in correcting pitch
66 halving and doubling errors, a large weight on the continuity constraint might introduce
67 errors^{12,27} by not recognizing gradual pitch transitions. Conversely, a weak continuity con-
68 straint may produce undesired fluctuations in the estimated pitch contour. These variations
69 in the pitch estimation using RAPT and PEFAC are illustrated in Figure 1. In box 4, the
70 estimated pitch from PEFAC is smoother than the ground truth, which has pitch transitions.
71 This could be due to the strong continuity constraint in the DP. However, in boxes 1 and
72 2 (RAPT), the inaccurate transitions could be due to relaxed continuity constraint. In box
73 3, the pitch estimation error occurs due to the absence of the pitch candidate. This could
74 be caused by inaccurate estimation or insufficient pitch candidates. Increasing the number
75 of pitch candidates would require a carefully designed cost function for the DP to result
76 in an accurate pitch contour. Hence, the effectiveness of a DP based approach depends
77 on the degree of the continuity constraint and the accuracy of pitch candidates and their
78 confidence-scores.

80 In this work, we propose a technique for computing pitch candidates and their confidence-
81 scores by combining complementary characteristics of two existing methods, namely, sub-
82 harmonic to harmonic ratio (SHR)¹⁰ and sawtooth wave inspired pitch estimator (SWIPE)⁹.

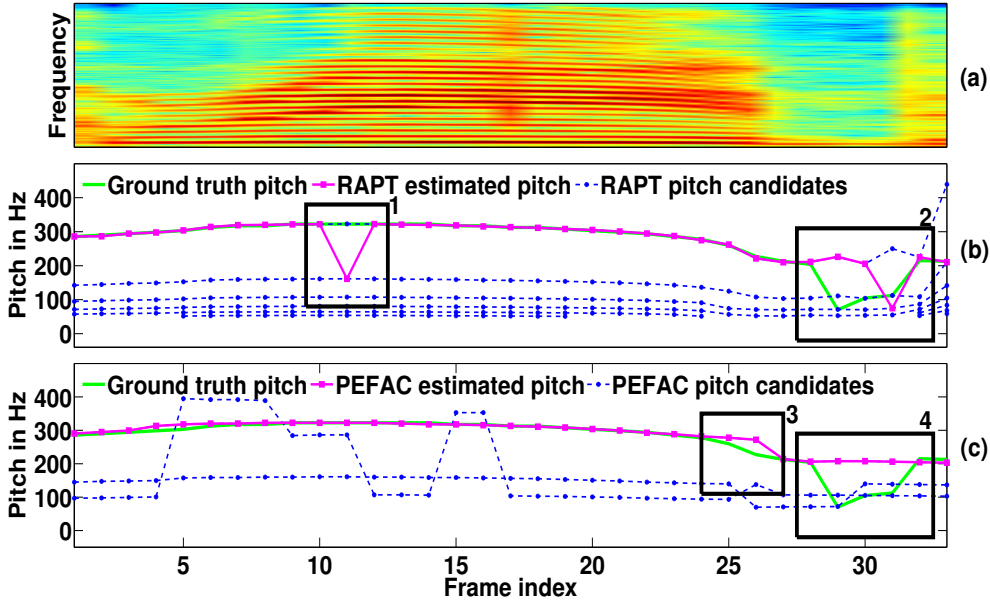


FIG. 1. An illustrative example describing two DP based techniques, namely RAPT and PEFAC –
a) spectrogram of an exemplary voiced segment using an FFT of 1024 with frame shift and length
of 10ms and 20ms respectively, b) pitch estimation by the RAPT algorithm, c) by the PEFAC
algorithm. The erroneous regions in the estimated pitch contour are indicated with black boxes
with box number at the top right corner.

83 Using these candidates, we employ a DP scheme to provide continuity in the pitch contour
84 only in a few selected frames, called DP frames, unlike a typical DP method that works for all
85 frames within a voiced segment. In the remaining non-DP frames, pitch is estimated using a
86 maximal confidence-score criterion. We observe that SHR achieves a significant accuracy in
87 the pitch estimation because it uses a good strategy for estimating reliable pitch candidates.
88 However, it only computes two candidates, causing estimation error in cases where the
89 ground truth pitch does not correspond to any of the candidates. We propose an extended
90 candidate estimation strategy based on SHR to increase the number of pitch candidates, such

91 that one of those candidates becomes more likely to correspond to the ground truth pitch.
92 Similarly, we extend the confidence-score computation strategy in SWIPE by exploiting the
93 window dependent properties (hanning window dependent kernel) and speech perception
94 and production based properties. The latter includes equivalent rectangular bandwidth
95 (ERB) frequency scale and decaying spectral envelope ($1/f$) similar to the glottal pulse
96 spectrum. These confidence-scores are also used to automatically determine the DP and
97 non-DP frames.

98 In addition to the proposed FSDP method for pitch estimation, we perform voiced-
99 unvoiced (VuV) classification in each frame using the pitch candidate confidence-scores.
100 Experiments for both pitch estimation and VuV classification are performed using three cor-
101 pora: KEELE²⁸, CSLU²⁹ and PaulBaghsaw (PB)³⁰ in clean as well as noisy conditions with
102 additive white Gaussian noise in 20, 10, 5 and 0dB SNRs. Gross pitch estimation (GPE)-20
103 error, root mean squared error (RMSE) and voiced and unvoiced (VuV) classification er-
104 ror are used as the evaluation metrics. We consider RAPT, PEFAC, SHR, and SWIPE as
105 the baseline schemes. For pitch estimation, the proposed FSDP is found to achieve lower
106 GPE-20 and RMSE compared to those of four baseline schemes, when the performance is
107 averaged across all SNR conditions. FSDP performs better than all four baseline schemes
108 for all three corpora in clean and in all SNR conditions, except for PaulBaghsaw corpus at
109 0dB SNR. For VuV classification, the proposed FSDP performs better than all four baseline
110 schemes for all three corpora in clean as well as all SNR conditions except at 20dB SNR on
111 CSLU corpus, where RAPT has the least VuV error.

112 II. PROPOSED FSDP APPROACH

113 The proposed FSDP approach has five stages, shown in Figure 2 and these stages are
 114 described using an exemplary voiced segment shown in Figure 3. The first stage computes
 115 pitch candidates $(p_t^k, 1 \leq k \leq K)$ at the t -th frame, where K is the total number of pitch
 116 candidates. In the second stage the confidence-score $C_t(k)$ associated to each candidate
 117 is computed. In the third stage, a VuV decision is taken at each frame based on the
 118 confidence-scores $C_t(k)$, and using a support vector machine (SVM) classifier, which was
 119 learnt in the training. This VuV decision is used in the fourth and fifth stages. We consider
 120 contiguous estimated voiced frames as one estimated voiced segment. Figure 3 shows the
 121 pitch candidates from the first stage for $K = 2$ in an estimated voiced segment. In the fourth
 122 stage, all frames in each estimated voiced segment are divided into two sets – DP frames
 123 and non-DP frames based on $C_t(k), 1 \leq k \leq K$. In Figure 3, the pitch candidates of the
 124 non-DP and DP frames are shown using red squares and blue diamonds respectively. The
 125 fifth stage estimates pitch (magenta line in Figure 3) for both types of frames separately.
 126 For the non-DP frames, pitch is estimated using the following maximization criteria:

$$k^{opt} = \arg \max_k C_t(k); \quad \hat{p}_t = p_t^{k^{opt}} \quad (1)$$

127 For the remaining frames, a DP based solution is used which selects one of the K pitch
 128 candidates in each frame such that the resultant pitch trajectory is maximally smooth within
 130 the segment.

131 It should be noted that the estimated unvoiced frames are not processed in the fourth
 132 and fifth stages. However, we use the maximization criteria in (1) to obtain pitch in the

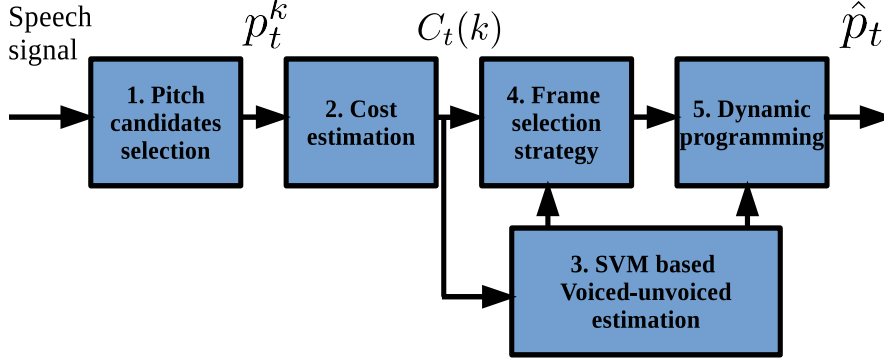


FIG. 2. Block diagram illustrating the steps of the FSDP method

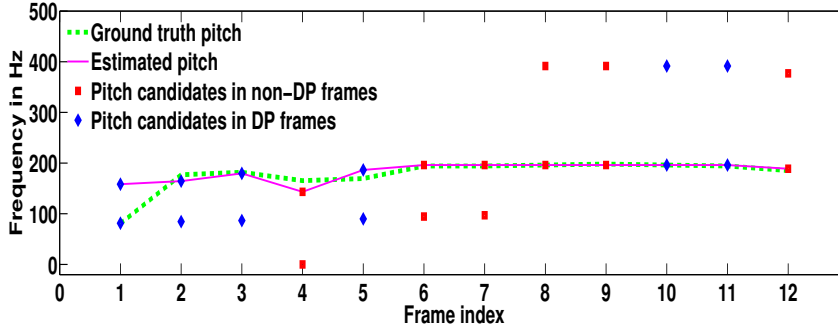


FIG. 3. An illustrative example explaining the proposed FSDP method

133 estimated unvoiced frames so that pitch is predicted in all frames of an utterance. This is
 134 done to obtain the pitch values at all ground truth voiced frames.

135 A. Pitch candidate selection

136 Pitch candidates are computed by following the two steps of the SHR method¹⁰. In the
 137 first step, we define $S_t(f)$ at the t -th frame as:

$$S_t(f) = \sum_{n=1}^N A_t(nf) - A_t\left(\left(n - \frac{1}{2}\right)f\right) \quad (2)$$

138 where $A_t(f)$ is the short time amplitude spectrum at the t -th frame and N is the maximum
139 number of harmonics contained in $A_t(f)$ ¹⁰. The $S_t(f)$ measures the difference between
140 amplitude sums at harmonic and at sub-harmonic components of the frequency f . This value
141 is expected to be maximum at the pitch frequency because a typical spectrum of a periodic
142 signal has high amplitudes at the harmonics of fundamental frequency and low amplitudes
143 at the sub-harmonics. In the case of non-periodic signals, for example an unvoiced sound,
144 the sum of the spectrum at the sub-harmonics would be relatively higher compared to that
145 of a periodic signal and, hence, the $S_t(f)$ might not be as high as that of voiced (periodic)
146 speech signal.

147 We observe that for some voiced speech segments the maximum of $S_t(f)$ may not cor-
148 respond to the pitch frequency. Hence, pitch estimation based on a strategy that selects
149 the frequency by maximizing $S_t(f)$ would introduce errors. We observe that most of these
150 errors are pitch halving and doubling, which are also common source of errors in most of
151 the existing pitch estimation methods¹⁹. This suggests that the candidate pitch frequency
152 could be obtained by multiplying frequency corresponding to the highest peak of $S_t(f)$ with
153 integer powers of 2.

154 In the second step, based on the above observation, we compute K pitch candidates as:

$$p_t^k = \begin{cases} \arg \max_f S_t(f) & \text{for } k = \lceil \frac{K}{2} \rceil \\ \arg \max_{\substack{k - \lceil \frac{K}{2} \rceil \leq f \leq k - \lfloor \frac{K}{2} \rfloor}} S_t(f) & \text{for } k \neq \lceil \frac{K}{2} \rceil \end{cases} \quad (3)$$

155 where p_t^k is k -th pitch candidate at t -th frame for $k \in \{1, 2, \dots, K\}$, $\lceil \frac{K}{2} \rceil$ is the small-
156 est integer greater than $\frac{K}{2}$ and $p_{t,Left}^{k - \lceil \frac{K}{2} \rceil}$ and $p_{t,Right}^{k - \lfloor \frac{K}{2} \rfloor}$ are equal to $(1 - \frac{1}{16}) = 0.9375$ and

157 $(1 + \frac{1}{16}) = 1.0625$ times $2^{k - \lceil \frac{K}{2} \rceil} p_t^{\lceil \frac{K}{2} \rceil}$ respectively. In particular, for $1 \leq k < \lceil \frac{K}{2} \rceil$, p_t^k in-
 158 cludes the frequencies around the sub-harmonics (negative integer powers of 2) of $p_t^{\lceil \frac{K}{2} \rceil}$ that
 159 fall within the frequency band ranging from $p_{t,Left}^{k - \lceil \frac{K}{2} \rceil}$ to $p_{t,Right}^{k - \lceil \frac{K}{2} \rceil}$, based on $\frac{1}{6}$ octave band at
 160 each candidate, which is a linear approximation for the critical bands of the ear^{31,32}. The
 161 $\frac{1}{6}$ octave band at $2^{k - \lceil \frac{K}{2} \rceil} p_t^{\lceil \frac{K}{2} \rceil}$ is equal to $(2^{\frac{1}{6}} - 1) \approx 0.125$ times $2^{k - \lceil \frac{K}{2} \rceil} p_t^{\lceil \frac{K}{2} \rceil}$, which is
 162 equal to $p_{t,Right}^{k - \lceil \frac{K}{2} \rceil} - p_{t,Left}^{k - \lceil \frac{K}{2} \rceil}$. Similarly, for $k > \lceil \frac{K}{2} \rceil$, p_t^k includes the frequencies around the
 163 harmonics (positive integer powers of 2) of $p_t^{\lceil \frac{K}{2} \rceil}$. We do not compute p_t^k beyond the typical
 164 pitch frequency range (50-550Hz). Hence the value of K is upper bounded by the total
 165 number of pitch candidates within the pitch range. The value of K is learnt in the training
 166 stage and is kept fixed for all the frames during the estimation of pitch and VuV decisions.

167 B. Candidate confidence-score computation

168 We modify the confidence-score computation steps in SWIPE⁹ and define the confidence-
 169 score ($C_t(k)$) associated with each pitch candidate as:

$$C_t(k) = \frac{\sum_{f'} \Phi(p_t^k, f') \sqrt{\Lambda_t(p_t^k, f')} \frac{1}{\sqrt{f'}}}{\left| \Phi^+(p_t^k, f') \frac{1}{\sqrt{f'}} \right| \left| \sqrt{\Lambda_t(p_t^k, f')} \right|} \quad (4)$$

170 where, $\Lambda_t(p_t^k, f')$ is the amplitude spectrum of a windowed speech signal at the t -th frame
 171 with frequency index f' . The amplitude spectrum ($\Lambda_t(p_t^k, f')$) is computed for every pitch
 172 candidate p_t^k using Hanning window of size equal to $\frac{8}{p_t^k}$. $\Phi^+(p_t^k, f')$ is the positive part of
 173 the kernel $\Phi(p_t^k, f')$ ⁹, which is defined for every pitch candidate p_t^k with frequency index f'

174 as $\sum_{i \in \{1\} \cup P} \Phi_i(p_t^k, f')$ where P is the set of prime numbers and $\Phi_i(p_t^k, f')$ defined as:

$$\Phi_i(p_t^k, f') = \begin{cases} \cos\left(2\pi \frac{f'}{p_t^k}\right), & \left|\frac{f'}{p_t^k} - i\right| < \frac{1}{4} \\ \frac{1}{2} \cos\left(2\pi \frac{f'}{p_t^k}\right), & \frac{1}{4} < \left|\frac{f'}{p_t^k} - i\right| < \frac{3}{4} \\ 0, & \textit{otherwise} \end{cases}$$

175 As demonstrated by Camacho et al⁹, the $C_t(k)$ is typically high at the pitch frequency when
 176 the window and the kernel are chosen appropriately.

177 C. VuV decisions estimation

178 We consider an SVM based classifier for estimating VuV decisions as a binary classification
 179 task. We obtain VuV decisions from the K candidate confidence-scores $C_t(k)$ belonging to
 180 each frame as a feature vector. Along with these feature vectors, we use ground truth
 181 VuV decisions labeled from the ground truth pitch values to train the SVM. In the labeling
 182 procedure, we consider the frames corresponding to zero pitch values as unvoiced and the
 183 remaining frames as voiced for all three corpora.

184 D. Frame selection strategy using nearest neighborhood

185 We observe that due to the mismatch between window and kernel choices, $C_t(k)$ could be
 186 high at a pitch candidate different from the correct pitch frequency. Thus, determining pitch
 187 frequency by finding the frequency corresponding to the highest confidence-score (SWIPE
 188 strategy) may not work uniformly well in all frames. We propose a method to automatically
 189 determine the frames (referred to as non-DP frames) where taking the frequency correspond-

190 ing to the highest confidence-score would accurately estimate the pitch frequency. In the
191 remaining frames (referred to as DP frames), we use DP for estimating pitch. For DP, confi-
192 dence scores are not used; rather, only pitch candidates are used. This helps in overcoming
193 the errors in the pitch estimated by SWIPE strategy in DP frames. Towards this, in the
194 training stage, we define two groups of pitch candidates – 1) the pitch candidate frequencies
195 lying within $\pm 20\%$ of the ground truth pitch called required pitch candidates (RPCs); 2)
196 other pitch candidates (non RPCs) away from (more than 20%) the ground truth pitch.
197 We refer to the voiced frames corresponding to RPCs with the highest confidence-score as
198 non-DP frames and the remaining voiced frames as DP frames and consider them as ground
199 truth DP and non-DP frames. In order to determine the DP and non-DP frames in testing
200 stage, we propose a frame selection strategy in the following section.

201 In the frame selection strategy, each frame of a voiced segment is categorized into either
202 a DP frame or a non-DP frame. For this, we utilize the confidence-score associated with
203 each candidate in developing the frame selection strategy. We use the confidence-scores
204 of all pitch candidates as K -dimensional feature vector and pose the frame selection as a
205 binary classification problem – non-DP frames as one class and DP frames as another class.
206 The classification is done using the nearest neighborhood (NN) classifier^{33,34} where r -nearest
207 neighbors are computed based on the Euclidean distance. The parameter r is learnt during
208 the training phase.

Algorithm 1 Pitch contour estimation algorithm based on DP

1: Initialization: $\mathcal{K} = \{1 : K\}$, $T = \text{length of voiced segment}$

2: **for** each voiced segment **do**

3: Initialization: $D_1(i) = 0 \quad \forall i \in \mathcal{K}$

4: **for** each frame t from 2 to T **do**

$$\forall i \in \mathcal{K}$$

5: **if** $t \in \text{DP frames}$ **then**

$$D_t(i) = \min_{j \in \mathcal{K}} \left\{ D_{t-1}(j) + \left(p_t^i - p_{t-1}^j \right)^2 \right\}$$
$$k_t(i) = \arg \min_{j \in \mathcal{K}} \left\{ D_{t-1}(j) + \left(p_t^i - p_{t-1}^j \right)^2 \right\}$$

6: **else**

$$k^{opt} = \arg \max_{j \in \mathcal{K}} \{ C_t(j) \}$$

$$p_t^i = p_t^{k^{opt}}$$

$$D_t(i) = \min_{j \in \mathcal{K}} \left\{ D_{t-1}(j) + \left(p_t^i - p_{t-1}^j \right)^2 \right\}$$
$$k_t(i) = \arg \min_{j \in \mathcal{K}} \left\{ D_{t-1}(j) + \left(p_t^i - p_{t-1}^j \right)^2 \right\}$$

7: **end if**

8: **end for**

9: Back tracking: $\eta_T = \arg \min_{i \in \mathcal{K}} \{ D_T(i) \}$, $\hat{p}_T = p_T^{\eta_T}$

10: **for** each frame t from $T - 1$ to 1 **do**

$$\eta_t = k_{t+1}(\eta_{t+1})$$

$$\hat{p}_t = p_t^{\eta_t}$$

11: **end for**

12: **end for**

209 E. Dynamic programming

210 Most of the pitch estimation algorithms are prone to octave errors, in which the estimated
211 pitch contour has abrupt transitions and differs from the original pitch by a factor of two or
212 a half^{12,19}. However, realistic pitch contour does not vary such abruptly and pitch variation
213 across frames is, in general, smooth in nature^{12,19,27}. In order to avoid these abrupt jumps
214 due to erroneous pitch estimates, we incorporate a temporal continuity constraint to estimate
215 the pitch in the DP frames. The continuity constraint is implemented using DP approach^{35,36}
216 with the Euclidean distance as an objective measure. The objective function involved in the
217 DP approach is given by

$$\hat{p}_t = \arg \min_{p_t; t \in \mathcal{F}} \sum_t (p_t - p_{t-1})^2 \quad (5)$$

such that $\hat{p}_t = p_t^{k^{opt}} \forall t \in \text{non-DP frames}$

218 where \mathcal{F} is a set of frames in a voiced region. The detailed algorithmic steps for solving (5)
219 are provided in Algorithm 1.

220 III. DATABASE

221 We use KEELE²⁸, CSLU²⁹, and PaulBaghsaw (PB)³⁰ corpora for all experiments in
222 this work. Table I shows the details of the three corpora and the number of recordings
223 considered in our experiments. In the experiments, we consider only the sentences belonging
224 to both the male and the female subjects from all three corpora and exclude the sentences
225 belonging to the children. In all three corpora, each spoken utterance has been recorded
226 simultaneously with a laryngograph signal, which is used to compute the reference pitch

227 considered as the ground truth. KEELE database consists of utterances from five male, five
 228 female and five children speakers reading “The North wind story”. CSLU database consists
 229 of 50 phonetically rich sentences spoken by seven male and five female speakers. These
 230 sentences have been collected from the TIMIT and Harvard Psychoacoustic corpora³⁷. Each
 231 speaker has uttered every sentence in three different contexts. PB database consists of 50
 232 sentences spoken by one male and one female speakers.

TABLE I. Details of the three corpora used in the experiments in this work

		KEELE	CSLU	PB
Number of sentences	Overall	15	1800	100
	considered	10	1800	100
Number of speakers	male	5	7	1
	female	5	5	1
	children	5	–	–
Availability of laryngograph		Yes	Yes	Yes

233 IV. EXPERIMENTAL RESULTS

234 A. Experimental setup

235 We compare the performance of the proposed FSDP for pitch estimation and VuV classi-
 236 fication with four existing methods (SHR, SWIPE, RAPT and PEFAC) using speech signal

237 in clean condition and in additive white Gaussian noise in four SNR conditions: 20, 10, 5
 238 and 0dB. KEELE, CSLU and PausBaghsaw (PB) corpora have been used for this purpose.
 239 Among the four existing methods, Matlab implementations of the four methods, namely,
 240 SHR, SWIPE, RAPT and PEFAC are directly available³⁸⁻⁴⁰ and are used for the compar-
 241 ison. The gross pitch error (GPE-20) and root mean squared error (RMSE)¹⁴ are used as
 242 the metrics for comparing the performance of pitch estimation using different methods. The
 243 GPE-20 is computed as $100 \times \frac{N_{err}}{N_v}$, where, N_{err} is the total number of erroneous frames, in
 244 which the estimated pitch values fall outside $\pm 20\%$ of the ground truth pitch value and N_v
 245 is the total number of voiced frames. Ground truth pitch is computed from the laryngo-
 246 graph signal available with individual corpus. Both GPE-20 and RMSE are computed by
 247 discarding the estimated pitch at the boundary frames (first and last frame) in every voiced
 248 segment. The parameters K and r are learnt using ground truth VuVs separately for each
 249 corpus from randomly chosen 20% data in clean condition among which 75% of the data is
 250 used for the training and the remaining used for the development. Among K and r , first, we
 251 obtain the best K which results in the least GPE-20 error on the entire 20% data. Then, the
 252 best r is learnt using the best K considering the errors computed on the development set.
 253 The parameters corresponding to the least GPE-20 error on a corpus are used to estimate
 254 pitch within the corpus and across several other corpora in clean and noisy conditions using
 255 estimated VuVs to examine the generalizability of the proposed method.

256 The performance of VuV classification using different methods are compared using the
 257 classification error⁴¹. We use SVM classifier with RBF kernel for the classification task with
 258 the complexity parameter (C) equal to 1.0 and with kernel coefficient (γ) equal to 1/number

259 of features. SVM classifier is implemented using Scikit-learn⁴². We train the SVM using a
260 training set identical to that for learning the parameters in the pitch estimation task. These
261 trained SVM models are used to estimate VuV decisions within the corpus and across the
262 corpora in clean and noisy conditions. During comparison, we use readily available VuV
263 decisions from all four existing methods except SWIPE for which classification error is not
264 reported.

265 The performance of the proposed FSDP method depends on the accuracies in the es-
266 timation of DP & non-DP frames and VuV decisions. To understand the effect of each
267 of these factors on the overall performance, we present the results in three sub-sections in
268 Section IV B. Section IV B 1 discusses the pitch estimation accuracy with ground truth DP
269 & non-DP frames and VuV decisions. Section IV B 2 discusses the effect of estimated DP &
270 non-DP frames on the overall performance. Similarly, Section IV B 3 explains the effect of
271 estimated DP & non-DP frames and VuV decisions. Following this, we analyze the reasons
272 for a better performance using the proposed FSDP methods over four baseline schemes in
273 two sub-sections – IV B 4 and IV B 5. For this analysis in Section IV B 4 and IV B 5, the
274 benefit of FSDP are highlighted by comparing with SHR & SWIPE and with RAPT &
275 PEFAC respectively. Note that, the performance of the proposed method also depends on
276 the accuracy of the pitch candidates and their confident scores, which is discussed in Section
277 IV B 6. Finally, in Section IV B 7, we present the accuracy of VuV classification.

TABLE II. GPE-20 obtained using the FSDP with ground truth DP and non-DP frames. A bold entry for a corpus and noise condition indicates the least GPE-20 among different K .

		$K = 2$	$K = 3$	$K = 4$
Clean	KEELE	0.79	0.74	0.75
	CSLU	1.16	0.91	0.95
	PB	1.34	1.33	1.25
20dB	KEELE	0.85	0.81	0.81
	CSLU	1.18	0.92	0.96
	PB	1.36	1.34	1.26
10dB	KEELE	1.20	1.18	1.20
	CSLU	1.49	1.14	1.20
	PB	1.90	1.89	1.78
5dB	KEELE	2.09	1.99	1.97
	CSLU	2.11	1.67	1.73
	PB	3.15	3.13	2.92
0dB	KEELE	5.29	5.08	4.86
	CSLU	3.88	3.16	3.14
	PB	6.25	6.21	5.56

TABLE III. GPE-20 obtained using the FSDP within and across all the three corpora using corpus specific parameters K and r learnt on the development set.

		FSDP		
		KEELE	CSLU	PB
clean	KEELE	0.79	1.04	1.17
	CSLU	1.61	1.52	1.74
	PB	1.49	1.45	1.36
20dB	KEELE	1.12	1.15	1.24
	CSLU	1.65	1.57	1.79
	PB	1.45	1.49	1.36
10dB	KEELE	1.56	1.60	1.67
	CSLU	2.01	2.02	2.15
	PB	2.03	2.02	1.92
5dB	KEELE	2.91	2.73	2.93
	CSLU	2.73	2.75	2.87
	PB	3.61	3.37	3.21
0dB	KEELE	6.90	6.48	6.59
	CSLU	4.78	4.73	4.88
	PB	7.57	6.88	6.42

B. Results and discussions

1. GPE-20 using FSDP with ground truth DP & non-DP frames and VuV deci-

sions

Frame selection is one of the key components in the proposed FSDP approach. An error in frame selection causes errors in pitch values estimated using FSDP. Hence, we first compute the GPE-20 using FSDP where we use the ground truth DP and non-DP labels and VuV decisions (i.e., no errors due to either automatic frame selection or VuV classification). This could be used as the lower bound on the GPE-20 of the FSDP scheme. Table II shows these GPE-20 values computed on entire data from three corpora under clean and all noisy conditions for $K \in \{n; 2 \leq n \leq 4\}$. It is clear from Table II that the least GPE-20 increases with decreasing SNR. It also varies across different corpora. From the table, it is observed that the best K (corresponding to the least GPE-20) is 3 in clean, 20dB and 10dB SNR conditions and 4 in 0dB SNR conditions for KEELE and CSLU corpora. For PB, the best K is found to be 4 in clean and all noisy conditions. This indicates that the best K varies even within a corpus in clean and all noisy conditions; it also varies across three corpora. However, we consider the best K obtained in clean condition for each corpus to find the best choice of parameter r for NN based frame selection strategy.

295 2. *GPE-20 using FSDP with estimated DP & non-DP frames and ground truth*

296 *VuV decisions*

297 The best choice of r is obtained for the frame selection strategy separately for each
298 corpus using ground truth VuV decisions. We find the best choice of the parameter
299 $r \in \{1 + 2n; 0 \leq n \leq 12\}$ on the development set for clean condition using GPE-20. The
300 parameters K and r corresponding to the minimum GPE-20 are found to be (3 and 1), (3
301 and 21) and (4 and 1) for KEELE, CSLU and PB corpora respectively. From the optimal
302 choice of r , it is observed that the parameter value changes in a corpus dependent manner.

303 Table III shows GPE-20 values on the entire data from all three corpora separately in
304 clean and noisy conditions using K and r learnt for each corpus. In the table, each column
305 indicates the corpus that is used for optimizing the parameters. It should be noted that
306 the parameters are optimized for clean conditions. The diagonal entries (shaded regions in
307 every 3×3 sub tables in Table III) indicate the GPE-20 values within the corpus (matched
308 development and test corpora) and the off-diagonal values indicate errors across corpora
309 (mis-matched development and test corpora). Bold entry for each corpus (every row) in
310 Table III indicates the least GPE-20 value among all columns, which indicates the best
311 development set. From the table, it is interesting to observe that the least errors are not
312 confined to diagonal entries only, particularly at low SNR.

TABLE IV. Comparison of pitch estimation and VuV classification performance of RAPT, PEFAC, SHR, SWIPE and FSDP. The performance of different methods is compared using GPE (%) and RMSE (%).

		KEELE			CSLU			PB		
		GPE-20	RMSE	VuV	GPE-20	RMSE	VuV	GPE-20	RMSE	VuV
Clean	RAPT	2.85	17.27	8.99	4.39	20.85	6.77	2.24	34.93	10.40
	PEFAC	12.16	41.21	12.68	5.93	27.99	10.12	3.39	15.79	9.32
	SHR	1.73	13.16	12.72	2.63	17.12	13.51	1.93	11.14	8.13
	SWIPE	4.31	21.69	–	3.41	22.34	–	2.52	15.95	–
	FSDP	0.89	8.90	6.43	1.65	13.31	5.91	1.50	9.01	7.29
20dB	RAPT	4.07	21.18	6.30	5.12	24.99	4.43	3.94	16.12	7.49
	PEFAC	12.35	46.58	12.81	6.09	28.26	10.31	3.39	16.17	9.25
	SHR	1.87	12.95	8.22	2.73	17.54	5.82	1.99	11.56	6.54
	SWIPE	4.68	21.94	–	3.78	22.63	–	2.83	16.99	–
	FSDP	1.17	9.91	6.12	1.69	13.54	5.45	1.48	8.85	6.51
10dB	RAPT	16.48	35.51	9.46	14.81	34.37	6.31	16.19	27.31	7.75
	PEFAC	11.91	41.13	13.79	6.56	28.28	11.19	3.72	16.02	10.02
	SHR	2.50	14.97	15.01	3.42	19.31	11.07	2.46	12.08	9.05
	SWIPE	8.12	28.56	–	6.02	26.44	–	4.42	21.02	–
	FSDP	1.59	11.39	7.09	2.05	14.56	5.58	2.00	9.60	5.87
5dB	RAPT	31.70	53.35	17.03	24.07	46.78	10.94	25.39	42.36	10.47
	PEFAC	12.69	39.47	15.00	7.24	28.10	12.11	4.58	17.22	11.01
	SHR	4.24	18.28	26.48	4.59	21.50	21.42	3.87	14.45	17.17
	SWIPE	15.07	39.10	–	10.72	33.62	–	8.09	28.42	–
	FSDP	2.79	13.43	9.73	2.79	16.01	6.82	3.39	12.12	6.69
0dB	RAPT	59.77	75.76	30.62	48.67	69.38	23.21	51.63	69.01	20.07
	PEFAC	14.65	37.73	16.90	8.53	28.28	13.41	6.38	19.28	12.00
	SHR	8.26	23.99	41.10	7.52	26.23	36.24	7.90	20.38	28.27
	SWIPE	30.68	55.55	–	23.33	48.80	–	20.99	45.80	–
	FSDP	6.56	18.81	16.26	4.77	19.15	10.61	6.91	17.18	10.87

313 **3. Comparison of GPE-20 and RMSE from FSDP and baseline schemes**

314 Once the corpus and the SNR for a given test utterance is known, an accurate pitch
315 contour could be achieved by using the parameters (K and r) corresponding to the least
316 GPE-20 values (marked in bold) in Table III. However these corpus dependent parameters
317 and the corresponding GPE-20 values might not be generalizable for unseen data. So, it

318 may not be fair to compare these corpus and SNR specific GPE-20 values with the GPE-20
 319 values computed using four baseline methods across all corpora and SNR conditions. Hence
 320 in FSDP, we consider one parameter set for frame selection strategy across all corpora and
 321 SNR conditions. This parameter set corresponds to the least GPE-20 value on the entire
 322 data among all corpora and SNR conditions (marked in blue in Table III). K and r in this
 323 parameter set are found to be the ones learnt on KEELE, i.e., $K=3$ and $r=1$. Using these
 324 parameters, we estimate VuV decisions using the SVM model learnt on KEELE with $K=3$.
 325 Following this, GPE-20 and RMSE are computed for all three corpora.

326 Table IV shows the GPE-20 and RMSE values obtained on the three corpora using the
 327 proposed FSDP and four baseline methods (RAPT, PEFAC, SHR and SWIPE) at various
 328 noisy levels and clean condition. In addition, we consider all frames as DP frames (i.e., no
 329 frame selection) in Equation 5 and compute the GPE-20 and RMSE to analyze the benefit
 330 of frame selection in FSDP scheme. However, pitch estimation using all frames as DP frames
 331 results in very poor performance; hence not reported in the table. The best performance
 332 for each metric is indicated in bold for each corpus and SNR condition. From the table,
 333 it is observed that the proposed FSDP performs better than baseline methods for all three
 334 corpora in clean and all SNR conditions except at 0dB SNR in PB corpus (GPE-20 value),
 335 at which PEFAC performs better than FSDP. When averaged across clean and all noisy
 336 conditions, FSDP achieves the least average GPE-20 and RMSE errors (2.60 and 12.49, 2.59
 337 and 15.31, 3.06 and 11.35) followed by SHR (3.72 and 16.67, 4.18 and 20.34, 3.63 and 13.92)
 338 for KEELE, CSLU and PB respectively. This implies that the strategies of SHR and SWIPE
 339 are complementary in nature and, when combined for computing pitch candidates and their

340 confidence-scores as in FSDP, they achieve better pitch estimation accuracy compared to
 341 the individual ones in most of the cases. The improvement in the performance of FSDP over
 342 the four baseline methods is analyzed separately in two following subsections by comparing
 343 with – 1) SHR and SWIPE (the variants of which have been used in FSDP) 2) RAPT and
 344 PEFAC (DP based methods).

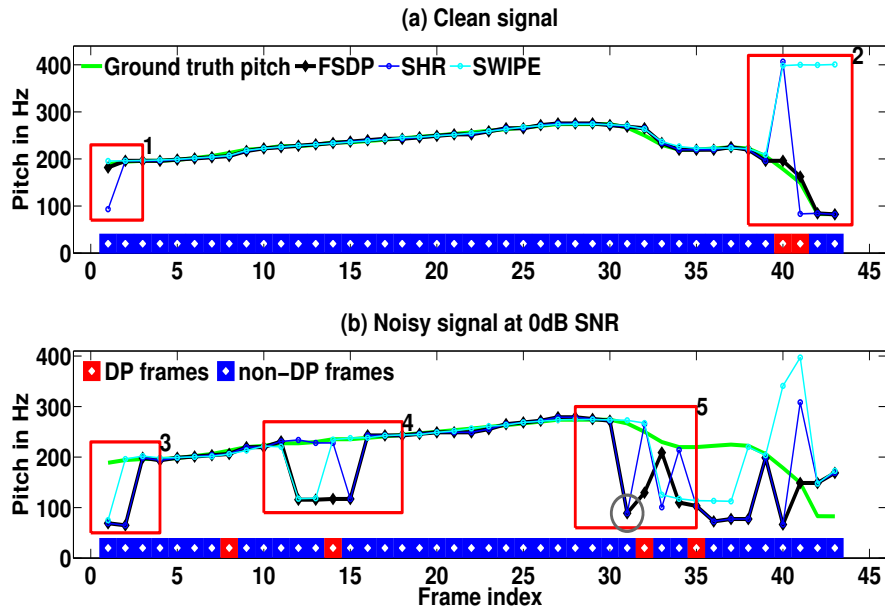


FIG. 4. Illustrative example describing the benefit of FSDP over SHR and SWIPE in a voiced segment. Red boxes 1, 2, 3, 4 and 5 are used to indicate the significant variations in the estimated pitch from the ground truth pitch. Red and blue horizontal patches indicate DP and non-DP frames respectively.

345 4. Comparison with SHR and SWIPE

346 Figure 4a and 4b show the estimated pitch trajectories for an exemplary voiced segment
 347 taken from the KEELE database in clean and noisy (SNR 0dB) conditions respectively. In

TABLE V. Comparison of erroneous frames (%) for both DP and non-DP categories obtained from FSDP in clean and all noisy conditions for all three corpora. All the percentages for each corpus are computed with respect to the total number of voiced frames

		estimated DP frames	EDPFs	ENDPFs	EDP_NDPFs (ENDPF's)
KEELE	clean	0.89	0.25	0.65	0.02 (0.01)
	20dB	0.98	0.25	0.92	0.02 (0.01)
	10dB	1.26	0.32	1.27	0.04 (0.03)
	5dB	2.08	0.46	2.33	0.12 (0.10)
	0dB	2.64	0.74	5.82	0.34 (0.28)
CSLU	clean	0.70	0.17	1.48	0.03 (0.02)
	20dB	0.77	0.17	1.52	0.03 (0.02)
	10dB	1.05	0.22	1.84	0.05 (0.04)
	5dB	1.46	0.26	2.53	0.05 (0.04)
	0dB	2.48	0.44	4.33	0.18 (0.13)
PB	clean	1.15	0.06	1.43	0.00 (0.00)
	20dB	1.29	0.09	1.38	0.00 (0.00)
	10dB	1.79	0.24	1.76	0.00 (0.00)
	5dB	1.92	0.41	2.98	0.03 (0.02)
	0dB	2.93	0.70	6.21	0.26 (0.22)

348 box-1 all methods estimate pitch correctly except the SHR. This indicates that original pitch
349 could be one of the pitch candidates in SHR, but the selection criteria used in SHR has led
350 to wrong estimation of the pitch. In box-2, where the ground truth pitch has large variation,
351 the proposed FSDP estimates pitch more accurately compared to all other methods. The
352 SWIPE estimates wrongly at most of the points, which could be due to the large amount of
353 errors in SWIPE when the actual pitch has wide variations. This could be because SWIPE
354 considers many pitch candidates for estimating the pitch. SHR has better pitch estimates
355 than those of SWIPE but worse than those of FSDP. When ground truth pitch has wide
356 variation, we observe that the estimates of the pitch candidates and their confidence-scores
357 become less reliable. This causes the SHR and SWIPE to result in octave errors. We

358 also observe that such unreliable frames often get classified as DP frames using the nearest
359 neighborhood strategy. Since the DP in the proposed scheme does not directly use the
360 confidence-scores of the pitch candidates in DP frames and rather uses estimated pitch from
361 neighboring non-DP frames to compensate the octave errors, the accuracy in the estimated
362 pitch in these unreliable frames improves by using FSDP.

363 From Figure 4b, it is observed that the estimation errors are more in 0dB SNR compared
364 to the clean condition for all the methods. This observation is consistent with the overall
365 performance degradation in Table IV from clean to 0dB SNR condition. The performance
366 degradation of FSDP could be due to two reasons. The first reason is that the estimated
367 DP frames are more (2 in Figure 4a and 4 in Figure 4b, as highlighted using red horizontal
368 patches) in case of 0dB SNR than in the clean condition. Higher number of DP frames could
369 cause a smooth pitch trajectory even in frames with large ground truth pitch variations, and
370 thereby resulting in a lower performance at 0dB SNR. The percentage of such DP frames
371 that cause errors in the pitch estimation, called erroneous DP frames (EDPFs), are listed in
372 the fourth column of Table V across all three corpora in clean and all noisy conditions. From
373 the table, it is observed that the DP frames and EDPFs increase from clean to 0dB SNR
374 condition for all three corpora. This implies that more DP frames result in more EDPFs,
375 and, hence, the performance could degrade from clean to 0dB SNR.

376 The second reason for poor performance of FSDP in low SNR condition could be a
377 large number of non-DP frames which result in pitch estimation errors at 0dB SNR, called
378 erroneous non-DP frames (ENDPFs), (0 in Figure 4a and 15 in Figure 4b). In the entire
379 set of ENDPFs, a subset of ENDPFs, indicated as ENDPF' (13-th, 15-th and 31-st frames

380 in Figure 4b), introduces pitch estimation errors in the neighboring DP frames due to the
 381 smoothing constraint in DP. For illustration, consider the 31-st frame marked in gray circle
 382 in Figure 4b in the box-5. This frame is classified as a non-DP frame (but it is ENDPF)
 383 by the nearest neighborhood frame selection strategy. Because of this, FSDP estimates
 384 the pitch values incorrectly at the neighboring DP frame (32-nd) in box-5 by following a
 385 wrong smooth trajectory. ENDPFs and ENDPF' are listed in the fifth and sixth column
 386 (in brackets) of Table V. The percentage of such DP frames that results in pitch estimation
 387 errors due to ENDPF's (indicated as EDP_NDPFs) are listed in the sixth column of Table
 388 V. From the table, it is observed that EDP_NDPFs are more than ENDPF's for all three
 389 corpora. This indicates that the number of frames with estimated erroneous pitch is more for
 390 every ENDPF' than that for every remaining ENDPFs. From the table, it is also observed
 391 that ENDPFs as well as ENDPF's gradually increase from clean to 0dB SNR for all three
 392 corpora. Hence the additional pitch estimation errors by the ENDPFs along with pitch
 393 estimation errors by the EDPFs could result in further performance degradation. These
 394 observations from EDPFs and ENDPFs are consistent with the performance degradation of
 395 FSDP in Table IV for all three corpora.

396 5. Comparison with RAPT and PEFAC

397 From Table IV, it is observed that the GPE-20 of RAPT varies largely from clean to
 398 0dB SNR condition compared to all other methods for all three corpora. This observation
 399 is consistent with the experimental findings by Gonzalez and Brookes¹². The worst perfor-
 400 mance of RAPT at 0dB SNR could be due to the increase in incorrect pitch candidates by

401 NCCF. PEFAC performs worse in the clean case but better in the noisy case compared to
 402 RAPT. This is because it was designed specifically for noisy signal with low SNRs. However,
 403 FSDP performs better in both clean and noisy conditions in almost all cases. This superior
 404 performance could be because FSDP performs DP only in the selected frames with few pitch
 405 candidates (optimal $K = 3$) using a few parameters (K and r).

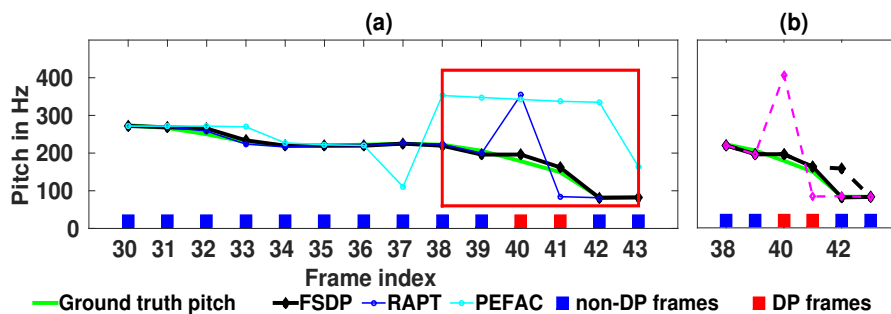


FIG. 5. Illustrative part of the voiced segment used in Figure 4 describing – a) the benefit of FSDP over RAPT and PEFAC b) the benefit of DP and non-DP frames in FSDP. The segment within red rectangular box in Figure 5a is shown in Figure 5b. The dotted black and magenta lines in Figure 5b indicate the estimated pitch trajectories when 40-th, 41-st and 42-nd frames are all DP frames and non-DP frames respectively.

406 Figure 5a shows the pitch trajectories obtained using RAPT, PEFAC and FSDP in clean
 407 condition for an exemplary voiced segment used in Figure 4. From Figure 5a, it is observed
 408 that PEFAC estimates incorrect pitch in the region, indicated by the red box in the figure,
 409 could be due to wrong pitch candidates in the highlighted region that result in a smooth
 410 trajectory away from the ground truth pitch. In the same region, RAPT estimates wrong
 411 pitch due to large deviations away from the ground truth pitch. However these errors are

412 compensated in FSDP by using two different strategies in DP and non-DP frames – DP
413 frames (40-th and 41-st) minimize the transitions and non-DP frames (38-th, 39-th, 42-
414 nd, 43-th) allow pitch transitions without any smoothness constraint. Thus, the proposed
415 FSDP allows pitch transitions as well as pitch smoothness in the right proportion using
416 frame selection strategy thereby achieving better pitch estimation accuracy.

417 We elaborate these benefits with the help of Figure 5b, where, in addition to the pitch
418 trajectory using FSDP, two other hypothetical trajectories (dotted-black and magenta) are
419 shown when 40-th, 41-st and 42-nd frames are all assumed to be DP frames and non-DP
420 frames respectively. It is clear that both trajectories suffer from pitch error either due to
421 smoothness constraint (in 42-nd frame when all are assumed to be DP frames) or due to
422 confidence-score maximization criterion (in 40-th and 41-st frames when all are assumed to
423 be non-DP frames). However, providing smoothness constraint only in selected DP frames
424 (as done in FSDP) results in an accurate pitch trajectory.

425 **6. *FSDP error analysis***

426 Overall, pitch estimation errors using FSDP depend on the strategies used in DP and non-
427 DP frames as well as the accuracy of the pitch candidates and their confidence-scores. We
428 categorize these errors into three types. – 1) Absence of RPCs as pitch candidate selection
429 strategy fails to detect them, 2) Estimated confidence-score associated with non-RPCs is the
430 highest among all candidates (even when RPCs are present) due to the errors in candidate
431 confidence-score estimation in the non-DP frames, 3) Selecting non-RPCs as the estimated

TABLE VI. Comparison of the number of GPE-20 frames belonging to three different types of errors occurred with different pitch candidates in clean and all noisy conditions for all three corpora.

		$K = 2$			$K = 3$			$K = 4$		
		Absence	with RPCs		Absence	with RPCs		Absence	with RPCs	
		of RPCs	in non-DP	in DP	of RPCs	in non-DP	in DP	of RPCs	in non-DP	in DP
KEELE	clean	0.41	0.23	0.32	0.26	0.39	0.24	0.26	0.39	0.32
	20dB	0.44	0.63	0.11	0.30	0.72	0.15	0.30	0.75	0.19
	10dB	0.81	0.63	0.18	0.67	0.71	0.21	0.67	0.81	0.21
	5dB	1.69	0.75	0.35	1.48	0.89	0.41	1.40	1.16	0.35
	0dB	4.73	1.23	0.55	4.35	1.51	0.71	4.06	1.84	0.68
CSLU	clean	0.87	0.75	0.04	0.57	0.98	0.09	0.57	1.12	0.05
	20dB	0.87	0.78	0.05	0.57	1.03	0.08	0.57	1.17	0.05
	10dB	1.14	0.89	0.06	0.78	1.18	0.09	0.77	1.32	0.06
	5dB	1.73	1.03	0.07	1.29	1.36	0.12	1.27	1.51	0.09
	0dB	3.40	1.36	0.12	2.75	1.81	0.19	2.69	2.03	0.17
PB	clean	1.25	0.20	0.00	1.20	0.26	0.00	1.11	0.25	0.01
	20dB	1.29	0.19	0.00	1.22	0.26	0.01	1.13	0.20	0.02
	10dB	1.83	0.20	0.00	1.77	0.25	0.00	1.64	0.25	0.04
	5dB	3.03	0.34	0.00	2.96	0.40	0.02	2.73	0.44	0.04
	0dB	5.97	0.89	0.03	5.85	1.01	0.04	5.26	1.09	0.11

432 pitch (even when RPCs are present) because of errors due to smoothing constraint using
 433 DP in DP frames.

434 Table VI shows the percentage of GPE-20 frames belonging to these three types of errors
 435 in FSDP for $K = 2, 3, 4$ in clean and all noisy conditions for all three corpora. From the
 436 table it is observed that the errors due to the absence of RPCs are significant in most of
 437 the cases (especially at 5dB and 0dB SNR conditions for all three corpora). The errors
 438 due to the absence of RPCs are crucial in the proposed FSDP, since they determine the
 439 pitch estimation errors when there is no error in both frame selection strategy and pitch
 440 estimation strategy at DP and non-DP frames. We investigate the reason for this error in
 441 detail with the help of Figure 6 with $K=2$ using $S_t(f)$ for two exemplary voiced frames from

442 the KEELE database. In the figure, the ground truth pitch frequency is indicated in green
 443 and the estimated pitch candidates are indicated in red. From Figure 6a, it is observed that
 444 the ground truth pitch frequency is closer to one of the pitch candidates. Hence, FSDP
 445 estimates the pitch accurately by choosing the correct pitch candidate. However in Figure
 446 6b, both pitch candidates are far from the ground truth pitch which implies that the pitch
 447 candidate selection fails to estimate the RPCs. Hence FSDP fails to estimate the correct
 448 pitch. This underlines the importance of the pitch candidate selection method.

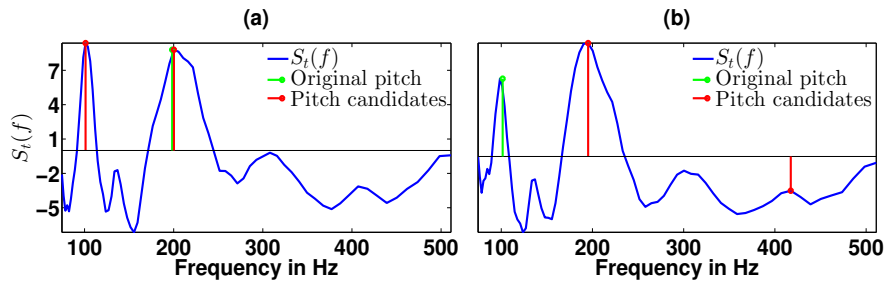


FIG. 6. Illustrative examples describing the importance of the pitch candidate selection method for $K = 2$ in which – a) the RPC is present b) the RPC is absent in the estimated pitch candidates

449 From Table VI, it is also observed that the errors due to the absence of RPCs reduce with
 450 increasing K for all SNRs and three corpora. This suggests that using more number of pitch
 451 candidates can reduce those errors. This could be because the search range of candidate
 452 selection method depends on K . Hence the RPCs which are missed with a low value of K
 453 can be detected with a high value of K . However, a high value of K does not guarantee a
 454 better pitch estimation accuracy due to increase in the second and third categories of error
 455 even when RPCs are present. This is also supported by the fact that the optimal K is found
 456 to be lower than 4 for all three corpora. Specifically, from Table VI, it is observed that

457 the second type of errors consistently increases with K in clean and all SNR conditions for
458 three corpora in almost all cases. This could be because the ambiguity in associating the
459 maximum confidence-score with the RPCs increases with increasing in K . Hence, the pitch
460 estimation method at non-DP frames fails to estimate the RPCs as a pitch.

461 The second and third categories of errors depend on the accuracy of the frame selec-
462 tion strategy and pitch estimation at DP and non-DP frames. We investigate these errors
463 in the proposed FSDP considering three different choices of DP and non-DP frames – 1)
464 ground truth DP and non-DP frames, 2) estimated DP and non-DP frames using the nearest
465 neighborhood 3) all frames as non-DP frames (to highlight the importance of DP frames).
466 We find the sum of second and third types of errors as 0.40, 1.01 and 1.04 respectively
467 for the above mentioned three choices, when averaged across all five SNRs and all three
468 corpora. Note that, considering all frames as DP frames results poor performance, hence
469 not reported. For computing these errors, the parameters of the proposed FSDP are kept
470 identical to those used in Table IV for all three corpora. It is observed that the average
471 errors increase (monotonically) from the first choice to the third choice. The non-zero error
472 in the first choice indicates that the errors are only due to incorrect pitch estimation at some
473 of DP and non-DP frames. A higher error in the second choice compared to the first choice
474 indicates combined effect of the errors caused by the frame classification strategy and pitch
475 estimation methods at DP and non-DP frames. Similarly, the highest error in the third
476 choice indicates that the pitch estimation errors due to the errors in the frame selection
477 strategy are less than those due to the pitch estimation strategy.

478 7. *VuV classification errors*

479 Table IV shows the VuV classification errors computed using FSDP and four baseline
480 methods for all three corpora. In the table, a bold entry for a given corpus and SNR
481 combination indicates the lowest VuV classification error. From the table, it is observed
482 that the proposed FSDP has the least VuV error in clean and all noisy conditions on all
483 three corpora except at 20dB SNR on CSLU corpus, where RAPT has the least VuV error.
484 This indicates that the proposed FSDP method performs better than the four baseline
485 schemes both in the pitch estimation and VuV classification tasks. It is interesting to notice
486 that no single baseline scheme has consistently performed the best among all the baseline
487 schemes across all noisy conditions on three corpora.

488 V. CONCLUSIONS

489 Realistic pitch trajectories are typically smooth in nature, but sometime they show large
490 variation in pitch values in a short span of time. In this work, we propose FSDP approach
491 for pitch estimation, which allows the estimated pitch trajectory to be smooth using DP
492 only in a few selected frames (called DP frames) unlike a typical DP based method which
493 forces the trajectory to be smooth over an entire voiced segment. In the remaining frames
494 (called non-DP frames), FSDP approach allows large variation in the estimated trajectory
495 by estimating pitch using a pitch candidate confidence-score maximization criterion where
496 the candidates and their confidence-scores are computed using variants of SHR and SWIPE.
497 These confidence-scores are used to automatically identify DP and non-DP frames. These

498 confidence-scores are also used for VuV classification using an SVM classifier. Experiments
499 with three corpora namely, KEELE, CSLU and PaulBaghsaw reveal that FSDP performs
500 better than four baseline methods considered in this work for pitch estimation as well as
501 VuV classification tasks.

502 The performance of the proposed FSDP method depends on the percentage of missing
503 RPCs, reliability in estimating the pitch candidate confidence-scores, classification accuracy
504 of DP and non-DP frames and effectiveness of smoothening constraint used in DP. Among
505 all the errors, the percentage of missing RPCs is found to be crucial, since these errors
506 determine the lower bound on the pitch estimation errors by the proposed method. Hence,
507 further investigation is required to reduce the missing RPCs with an appropriate candidate
508 selection strategy. In addition to this, the frame selection strategy needs to be improved.
509 Most of the errors in the nearest neighborhood based frame selection strategy is due to the
510 misclassification of DP and non-DP frames. Also, the computational cost involved in the
511 frame selection strategy is quite high. This is because the nearest neighborhood classifier in
512 the frame selection strategy computes a distance for each frame with all training samples.
513 The training set for the DP and non-DP frame classification is also found to be imbalanced
514 with a ratio of 1:100. Hence, a noise robust classifier with less computational complexity
515 under large imbalanced training set would be effective.

516 **VI. REFERENCES**

517 **REFERENCES**

518 ¹M. Wang and M. Lin, “An analysis of pitch in chinese spontaneous speech,” International
519 Symposium on Tonal Aspects of Languages: With Emphasis on Tone Languages (2004).

520 ²I. R. Murray and J. L. Arnott, “Toward the simulation of emotion in synthetic speech: A
521 review of the literature on human vocal emotion,” The Journal of the Acoustical Society
522 of America **93**(2), 1097–1108 (1993).

523 ³K. De Bot, “Visual feedback of intonation I: Effectiveness and induced practice behavior,”
524 Language and Speech **26**(4), 331–350 (1983).

525 ⁴A. Askenfelt, “Automatic notation of played music: the VISA project,” *Fontes Artis*
526 *Musicae* 109–120 (1979).

527 ⁵R. B. Dannenberg, W. P. Birmingham, G. Tzanetakis, C. Meek, N. Hu, and B. Pardo,
528 “The musart testbed for query-by-humming evaluation,” *Computer Music Journal* **28**(2),
529 34–48 (2004).

530 ⁶E. Yumoto, W. J. Gould, and T. Baer, “Harmonics-to-noise ratio as an index of the degree
531 of hoarseness,” *The journal of the Acoustical Society of America* **71**(6), 1544–1550 (1982).

532 ⁷G. Fant, *Acoustic theory of speech production: with calculations based on X-ray studies of*
533 *Russian articulations*, Vol. 2 (Walter de Gruyter, 1971).

534 ⁸A. Camacho, “SWIPE: A sawtooth waveform inspired pitch estimator,” Ph.D. thesis,
535 University of Florida (2007).

- 536 ⁹A. Camacho and J. G. Harris, “A sawtooth waveform inspired pitch estimator for speech
537 and music,” *The Journal of the Acoustical Society of America* **124**(3), 1638–1652 (2008).
- 538 ¹⁰X. Sun, “Pitch determination and voice quality analysis using subharmonic-to-harmonic
539 ratio,” *IEEE International Conference on Acoustics, Speech, and Signal Processing* **1**,
540 333–336 (2002).
- 541 ¹¹A. De Cheveigné and H. Kawahara, “YIN, a fundamental frequency estimator for speech
542 and music,” *The Journal of the Acoustical Society of America* **111**(4), 1917–1930 (2002).
- 543 ¹²S. Gonzalez and M. Brookes, “PEFAC—a pitch estimation algorithm robust to high levels
544 of noise,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing* **22**(2),
545 518–530 (2014).
- 546 ¹³O. Deshmukh, C. Y. Espy-Wilson, A. Salomon, and J. Singh, “Use of temporal information:
547 Detection of periodicity, aperiodicity, and pitch in speech,” *IEEE Transactions on Speech
548 and Audio Processing* **13**(5), 776–786 (2005).
- 549 ¹⁴C. Shahnaz, W.-P. Zhu, and M. O. Ahmad, “Pitch estimation based on a harmonic si-
550 nusoidal autocorrelation model and a time-domain matching scheme,” *IEEE Transactions
551 on Audio, Speech, and Language Processing* **20**(1), 322–335 (2012).
- 552 ¹⁵W. Bauer and W. Blankenship, “Dyptrack—a noise-tolerant pitch tracker,” Dept. of Defence
553 (NSA), Washington, USA, Unclassified Rep. NASL-S-210 (1974).
- 554 ¹⁶H. Ney, “A dynamic programming technique for nonlinear smoothing,” *IEEE International
555 Conference on Acoustics, Speech, and Signal Processing* **6**, 62–65 (1981).

- 556 ¹⁷H. Ney, “Dynamic programming algorithm for optimal estimation of speech parameter
557 contours,” *IEEE Transactions on Systems, Man and Cybernetics*, **2**, 208–214 (1983).
- 558 ¹⁸L. Sukhostat and Y. Imamverdiyev, “A comparative analysis of pitch detection methods
559 under the influence of different noise conditions,” *Journal of Voice* **29**(4), 410–417 (2015).
- 560 ¹⁹M. Asgari and I. Shafran, “Improving the accuracy and the robustness of harmonic model
561 for pitch estimation.,” *Proceedings Interspeech 1936–1940* (2013).
- 562 ²⁰D. Talkin, *A robust algorithm for pitch tracking (RAPT)* (W.B. Kleijin and K. K. Paliwal
563 Eds. Amsterdam, The Netherlands: Elsevier, 1995), pp. 495–518.
- 564 ²¹E. Azarov, M. Vashkevich, and A. Petrovsky, “Instantaneous pitch estimation based
565 on RAPT framework,” *European Signal Processing Conference (EUSIPCO)* 2787–2791
566 (2012).
- 567 ²²K. Han and D. Wang, “Neural network based pitch tracking in very noisy speech,”
568 *IEEE/ACM Transactions on Audio, Speech, and Language Processing* **22**(12), 2158–2168
569 (2014).
- 570 ²³H. Su, H. Zhang, X. Zhang, and G. Gao, “Convolutional neural network for robust pitch
571 determination,” *IEEE International Conference on Acoustics, Speech, and Signal Process-*
572 *ing* 579–583 (2016).
- 573 ²⁴K. Han and D. Wang, “Neural networks for supervised pitch tracking in noise,” *IEEE*
574 *International Conference on Acoustics, Speech, and Signal Processing* 1488–1492 (2014).
- 575 ²⁵K. Kasi and S. A. Zahorian, “Yet another algorithm for pitch tracking,” *IEEE International*
576 *Conference on Acoustics, Speech, and Signal Processing* **1**, 361–364 (2002).

- 577 ²⁶H. Ba, N. Yang, I. Demirkol, and W. Heinzelman, “BaNa: A hybrid approach for noise
578 resilient pitch detection,” IEEE Statistical Signal Processing Workshop (SSP) 369–372
579 (2012).
- 580 ²⁷L. Dolansky and P. Tjernlund, “On certain irregularities of voiced-speech waveforms,”
581 IEEE Transactions on Audio and Electroacoustics, **16**(1), 51–56 (1968).
- 582 ²⁸F. Plante, M. G, and A. W.A, “A pitch extraction reference database,” Proceedings Eu-
583 rospeech 95 837–840 (1995).
- 584 ²⁹A. Kain, “CSLU: Voices,” Linguistic Data Consortium, Philedelphia, PA, USA (2006).
- 585 ³⁰P. C. Bagshaw, S. M. Hiller, and M. A. Jack, “Enhanced pitch tracking and the processing
586 of f0 contours for computer aided intonation teaching.,” Proceedings of the Third European
587 Conference on SpeechCommunications and Technology 1003–1006 (1993).
- 588 ³¹F. A. Everest, K. C. Pohlmann, and T. Books, *The master handbook of acoustics*, Vol. 4
589 (McGraw-Hill New York, 2001).
- 590 ³²H. Fletcher, *Speech and hearing in communication* . (D. van Nostrand, 1953).
- 591 ³³T. M. Cover and P. E. Hart, “Nearest neighbor pattern classification,” Information Theory,
592 IEEE Transactions on **13**(1), 21–27 (1967).
- 593 ³⁴J. H. Friedman, J. L. Bentley, and R. A. Finkel, “An algorithm for finding best matches in
594 logarithmic expected time,” ACM Transactions on Mathematical Software (TOMS) **3**(3),
595 209–226 (1977).
- 596 ³⁵R. E. Bellman and S. E. Dreyfus, *Applied dynamic programming* (Rand Corporation,
597 1962).

- 598 ³⁶R. Bellman, *Dynamic programming* (Dover publications, 1957).
- 599 ³⁷A. B. Kain, “High resolution voice transformation,” Oregon Health & Science University
600 (2001).
- 601 ³⁸X. Sun, “Pitch determination algorithm,” Software, available [Jan 2016] from
602 <http://in.mathworks.com/matlabcentral/fileexchange/1230-pitch-determination->
603 [algorithm/content/shrp.m](http://in.mathworks.com/matlabcentral/fileexchange/1230-pitch-determination-algorithm/content/shrp.m) .
- 604 ³⁹A. Camacho, “SWIPE pitch estimation algorithm,” Software, available [Jan 2016] from
605 <http://www.cise.ufl.edu/acamacho/publications/swipep.m> .
- 606 ⁴⁰M. Brookes, “VOICEBOX: A speech processing toolbox for matlab. 2006,” URL
607 <http://www.ee.ic.ac.uk/hp/staff/dmb/voicebox/voicebox.html>. Available online (2003).
- 608 ⁴¹T. Drugman and A. Alwan, “Joint robust voicing detection and pitch estimation based on
609 residual harmonics.,” *Interspeech 1973–1976* (2011).
- 610 ⁴²F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
611 P. Prettenhofer, R. Weiss, and V. Dubourg, “Scikit-learn: Machine learning in python,”
612 *Journal of Machine Learning Research* **12**, 2825–2830 (2011).