# MATCHING PURSUIT FOR CHANNEL SELECTION IN COCHLEAR IMPLANTS BASED ON AN INTELLIGIBILITY METRIC

*C. H. Taal\**

Royal Institute of Technology (KTH)
Sound and Image Processing Lab
Stockholm, Sweden

*R. C. Hendriks and R. Heusdens*

Delft University of Technology
Signal and Information Processing Lab
Delft, the Netherlands

## ABSTRACT

In this paper the earlier proposed short-time objective intelligibility predictor (STOI) is simplified such that it can be expressed as a weighted $\ell_2$ norm in the auditory domain. Due to the mathematical properties of a norm, STOI can now be used with the matching pursuit algorithm in the *n-of-m* channel selection technique as found in several cochlear implant (CI) coding strategies. With this technique only a subset of frequency channels (electrodes) are stimulated, such that important channels can be updated more frequently and less significant channels are omitted. Intelligibility predictions with acoustic CI-simulations for normal-hearing listeners indicate that more intelligible speech is obtained with the proposed method compared to a conventional channel selection method based on peak picking. Reasons for this difference in performance are: (1) STOI considers an analysis window of a few hundreds of milliseconds in order to account for important low temporal modulations for speech intelligibility and (2) spectral leakage per channel is accounted for in the mathematical optimization process.

*Index Terms—* Speech intelligibility metric, matching pursuit, cochlear implants, channel selection.

## 1. INTRODUCTION

Reliable machine-driven predictors of speech intelligibly are of great interest in the design process of new speech processing algorithms, e.g., as used in mobile telephony, hearing aids or cochlear implants (CIs). They might replace costly and time consuming listening tests, at least in some stages of the algorithm development process. The drawback of many intelligibility predictors is that they are complex [1, 2] and do not have certain (mathematical) properties in order to derive optimal signal processing solutions, e.g., least-squares solutions. In previous work we proposed a short-time objective intelligibility (STOI) measure which can accurately predict the effect of background noise and various (non-)linear speech processing algorithms on speech intelligibility [3]. We will show

that STOI can be simplified to a weighted $\ell_2$ norm in the auditory domain which makes the measure mathematically tractable. Since STOI shows high correlation with the intelligibility of vocoded speech [3], as typically used in acoustic CI-simulations, the norm will be applied in the channel-selection technique with CI simulations [4, 5].

The channel-selection technique is also referred to as the *n-of-m* strategy where $n$ channels of the available $m$ frequency channels (electrodes) are stimulated, such that important channels can be updated more frequently and less significant channels are omitted. Different strategies exist to select those channels, e.g., based on peak-picking [6], psychoacoustic models [7] and other techniques [4]. However, those techniques optimize for certain (psychoacoustic) criteria which exclude important properties relevant for speech intelligibility [8]. For example, criteria relevant for speech intelligibility should take into account temporal modulation frequencies important for intelligibility (4-32 Hz) [9] and correlation based comparisons should be used rather than comparisons based on squared errors [8]. The proposed norm based on STOI takes into account these aspects.

Due to the mathematical properties of a norm, the channel selection can now be solved in an optimal manner for STOI with the matching pursuit algorithm [10]. Within this framework the electrical spread per electrode can also be easily taken into account, which is typically not part of the optimization process in existing *n-of-m* strategies. It will be shown that the proposed method leads to more intelligible speech compared to a general peak-picking algorithm by means of acoustical CI-simulations with normal-hearing listeners.

## 2. DERIVATION OF INTELLIGIBILITY METRIC

We will first introduce a general notation and explain the auditory model as used in STOI. Let $x(n)$ and $y(n)$ denote a clean and degraded speech signal, respectively, with time-sample index $n$, where $y$ is a vocoded version of $x$. A basic auditory model is applied to both signals in order to obtain an internal representation. Here, we only explain the notation for the internal representation of $x$. Similar definitions hold

for $y$. Let $\hat{x}_m(k)$ denote the $k^{th}$ DFT-bin of $xw_m$, where $w_m$ denotes a Hann-window function with frame-index $m$. Here, a frame length of 16 ms is used with 50% overlap. The short-time DFT spectrum is converted into auditory bands as follows:

$$X_{i,m} = \sum_k \left| \hat{h}_i(k)\hat{x}_m(k) \right|^2, \qquad (1)$$

where $i$ denotes the auditory band index and $\hat{h}_i$ represents an approximation of the magnitude response of a $4^{th}$ order gammatone filter as described in [11]. The value $X_{i,m}$ will be referred to as a time-frequency (TF) unit. In total, 32 filters are used with center frequencies linearly spaced on an ERB scale between 150 and 5000 Hz. STOI compares the clean and degraded speech in the auditory domain in blocks of approximately 400 milliseconds (see next section for more details). The following vector notation is used to denote such a block within one auditory band,

$$\mathbf{x}_{i,m} = \begin{bmatrix} X_{i,m-M+1} & X_{i,m-M+2} & \cdots & X_{i,m} \end{bmatrix}^T, \qquad (2)$$

where $M$ can be used to control the length of such a speech segment, depending on the sample rate and window size. In this work, a sample rate of 16 kHz is used where $M = 48$. Vectors are concatenated over all auditory bands to denote a complete TF-block as:

$$\mathbf{x}_m = \begin{bmatrix} \mathbf{x}_{1,m}^T & \mathbf{x}_{2,m}^T & \cdots & \mathbf{x}_{I,m}^T \end{bmatrix}^T, \qquad (3)$$

where $I = 32$ denotes the total amount of auditory filters. The operator notation $\mathbf{x}_m = \mathcal{I}_m\{x\}$ is used to denote the complete transform from the time-domain to one TF-block in the auditory domain.

## 2.1. STOI Background and Simplification

As proposed in STOI [3], an intermediate measure relevant for speech intelligibility of one TF-unit is defined as the sample correlation coefficient between the clean $(\mathbf{x}_{i,m})$ and degraded $(\mathbf{y}_{i,m})$ speech temporal band envelopes in one block. Blocks of a few hundreds of milliseconds are used to include important modulation frequencies for intelligibility [9]. The correlation coefficient is used, rather than, e.g., a squared error, to make sure that the measure is insensitive to band-level differences between $x$ and $y$, which should not have a strong impact on speech intelligibility [8]. To simplify, the correlation coefficient is defined on the magnitude squared envelopes rather than the magnitude envelopes, as was originally proposed in STOI [3]. The benefit of this choice will become clear in Section 4. This gives:

$$\rho_{i,m}(x,y) = \frac{\left\langle \mathbf{x}_{i,m} - \mu_{\mathbf{x}_{i,m}}, \mathbf{y}_{i,m} - \mu_{\mathbf{y}_{i,m}} \right\rangle}{\sigma_{\mathbf{x}_{i,m}} \sigma_{\mathbf{y}_{i,m}}}, \qquad (4)$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product with $\|\cdot\|$ as its induced $\ell_2$-norm, $\mu_{\mathbf{x}_{i,m}}$ the sample mean of $\mathbf{x}_{i,m}$ and $\sigma_{\mathbf{x}_{i,m}} = \left\| \mathbf{x}_{i,m} - \mu_{\mathbf{x}_{i,m}} \right\|$. Similar definitions hold for the degraded speech. The correlation coefficients $\rho_{i,m}(x,y)$ are then combined into one number by computing its average over all TF-units:

$$D = \frac{1}{\mathcal{M}} \sum_{i,m} \rho_{i,m}(x,y), \qquad (5)$$

where $\mathcal{M}$ denotes the total number of TF-blocks. It is expected that $D$ is a monotonically increasing function of the speech intelligibility of $y$. In computing $D$ only those TF-blocks are considered in the summation where speech is present, see [3] for more details. An additional clipping procedure in STOI, which was included to limit the intermediate intelligibility range, is discarded in this work for simplicity.

## 2.2. Interpretation as weighted $\ell_2$ norm

To rewrite the intelligibility measure as a norm we first express (4) as an inner product:

$$\rho_{i,m}(x,y) = \left\langle \bar{\mathbf{x}}_{i,m}, \bar{\mathbf{y}}_{i,m} \right\rangle, \qquad (6)$$

where a general normalization procedure is denoted by $\bar{(\cdot)} = \left( (\cdot) - \mu_{(\cdot)} \right) / \sigma_{(\cdot)}$. Hence, the inner product $\left\langle \bar{\mathbf{x}}_{i,m}, \bar{\mathbf{y}}_{i,m} \right\rangle$ can be used to induce the following norm:

$$\begin{aligned} \left\| \bar{\mathbf{x}}_{i,m} - \bar{\mathbf{y}}_{i,m} \right\|^2 &= \left\| \bar{\mathbf{x}}_{i,m} \right\|^2 + \left\| \bar{\mathbf{y}}_{i,m} \right\|^2 - 2\left\langle \bar{\mathbf{x}}_{i,m}, \bar{\mathbf{y}}_{i,m} \right\rangle \\ &= 2 - 2\rho_{i,m}(x,y). \end{aligned} \qquad (7)$$

It can now be observed that maximizing $\rho_{i,m}$ implies minimizing the norm $\left\| \bar{\mathbf{x}}_{i,m} - \bar{\mathbf{y}}_{i,m} \right\|^2$. However, its minimizing argument only determines the optimal $\mathbf{y}_{i,m}$ up to a scaling $\sigma_{\mathbf{y}_{i,m}}$ and amplitude shift $\mu_{\mathbf{y}_{i,m}}$. In this work we aim for the solution where the clean speech is the target, with the assumption that $\mu_{\mathbf{x}_{i,m}} \approx \mu_{\mathbf{y}_{i,m}}$ and $\sigma_{\mathbf{x}_{i,m}} \approx \sigma_{\mathbf{y}_{i,m}}$. This is motivated by the fact that we are working in blocks of a few hundreds of milliseconds, and it is expected that the errors introduced to $\mathbf{y}_{i,m}$ will average to a minimal impact when summing over all its elements in the calculation of the scaling $\sigma_{\mathbf{y}_{i,m}}$ and amplitude shift $\mu_{\mathbf{y}_{i,m}}$. This gives:

$$\left\| \bar{\mathbf{x}}_{i,m} - \bar{\mathbf{y}}_{i,m} \right\|^2 \approx \left\| a_{i,m}\left( \mathbf{x}_{i,m} - \mathbf{y}_{i,m} \right) \right\|^2 \qquad (8)$$

where $a_{i,m} = \sigma_{\mathbf{x}_{i,m}}^{-1}$. By vector concatenation as in (3) the summation over frequency $i$ in (5) can be replaced by defining a new norm over a complete TF-block. First, a diagonal weighting matrix is defined as:

$$\mathbf{A}_m = \mathrm{diag}\begin{pmatrix} a_{1,m}\mathbf{I}_M & a_{2,m}\mathbf{I}_M & \cdots & a_{I,m}\mathbf{I}_M \end{pmatrix}, \quad (9)$$

where $\mathbf{I}_M$ is the identity matrix of size $M$. A weighted norm for one TF-block is then given as follows:

$$\|\mathbf{A}_m\left(\mathbf{x}_m - \mathbf{y}_m\right)\|^2 = \sum_i \|a_{i,m}\left(\mathbf{x}_{i,m} - \mathbf{y}_{i,m}\right)\|^2. \quad (10)$$

These weighted norms are then combined by a summation over time, where for optimization purposes the averaging constant $\mathcal{M}$ in (5) can be discarded. Note that $\mathbf{A}_m$ is only a function of the clean speech $\mathbf{x}_m$. As a result, it only has to be calculated once for each frame after which the norm can be evaluated for any arbitrary $\mathbf{y}_m$.

## 3. APPLICATION TO CI CHANNEL SELECTION

The proposed intelligibility metric will be used in the CI channel-selection technique with the matching pursuit algorithm [10]. With this algorithm, a signal $x$ is synthesized as a weighted sum of functions (sometimes called atoms or elements) which are chosen from a dictionary [10]. The algorithm is iterative, where for each iteration $p$ the best matching function $g$ from the dictionary $\mathcal{D}$ is chosen and subtracted from the residual at the previous iteration. Since only one element is considered per iteration, the algorithm is greedy. The eventual synthesized speech signal can be described by:

$$x \approx \sum_p \alpha^{(p)} g^{(p)}, \quad (11)$$

where the selection of the best dictionary element and weighting coefficient $\alpha$ is based on minimizing some norm of the eventual residual $r$. For the $(p+1)^{th}$ iteration this residual is given as follows:

$$r^{(p+1)} = r^{(p)} - \alpha^{(p)} g^{(p)}, \quad (12)$$

where for the first iteration the residual is taken equal to the target signal, i.e, $r^{(1)} = x$. The optimal solution for the weighting coefficient and selection of the dictionary element in each iteration is given by [10],

$$\begin{aligned} \alpha^{(p)} &= \frac{\langle g^{(p)}, r^{(p)} \rangle}{\|g^{(p)}\|^2} \\ g^{(p)} &= \underset{g \in \mathcal{D}}{\arg\max} \frac{|\langle g, r^{(p)} \rangle|}{\|g\|} \end{aligned} \quad (13)$$

### 3.1. Intelligibility Relevant Matching Pursuit

Since all diagonal elements of $\mathbf{A}_m$ in (10) are real and positive a new norm relevant for speech intelligibility can be defined, say $\|\cdot\|_{\mathbf{A}_m}$, which is induced from the following inner product:

$$\langle \mathbf{x}_m, \mathbf{y}_m \rangle_{\mathbf{A}_m} = \langle \mathbf{A}_m \mathbf{x}_m, \mathbf{A}_m \mathbf{y}_m \rangle. \quad (14)$$

Now we can insert the proposed norm and inner product based on STOI in (12) and (13). Here, the dictionary will be defined by $\mathcal{D} = \mathbf{g}(\gamma)_{\gamma \in \Gamma}$, where $\Gamma$ denotes the set of CI frequency channel indices. Each element represents the internal representation of a short-time pulse within a specific CI channel and will be used to model $\mathbf{x}_m$. One can choose the dictionary according to the properties of the CI and include aspects like the pulse duration, channel center frequencies or the amount of current spread. To imply low algorithmic delay no future time-samples are taken into account in these internal representations for a given pulse.

For the first iteration where no channel selection has been made yet, the residual is set to $\mathbf{r}_m^{(1)} = \mathbf{x}_m$, where for the next iterations we have:

$$\mathbf{r}_m^{(p+1)} = \mathbf{r}_m^{(p)} - \alpha^{(p)} \mathbf{g}^{(p)}. \quad (15)$$

The solution for the best dictionary element and optimal weighting for each iteration relevant for the proposed metric is then given by:

$$\begin{aligned} \alpha^{(p)} &= \frac{\langle \mathbf{g}^{(p)}, \mathbf{r}_m^{(p)} \rangle_{\mathbf{A}_m}}{\|\mathbf{g}^{(p)}\|_{\mathbf{A}_m}^2} \\ \mathbf{g}^{(p)} &= \underset{g \in \mathcal{D}}{\arg\max} \frac{|\langle \mathbf{g}, \mathbf{r}_m^{(p)} \rangle_{\mathbf{A}_m}|}{\|\mathbf{g}\|_{\mathbf{A}_m}}. \end{aligned} \quad (16)$$

After the channels have been selected, the eventual residual $\mathbf{r}_m$ is stored and shifted one time-frame over $m$ for the initial residual $\mathbf{r}_{m+1}^{(1)}$. In this manner, past channel selections are also taken into account for the decisions of the current time-frame.
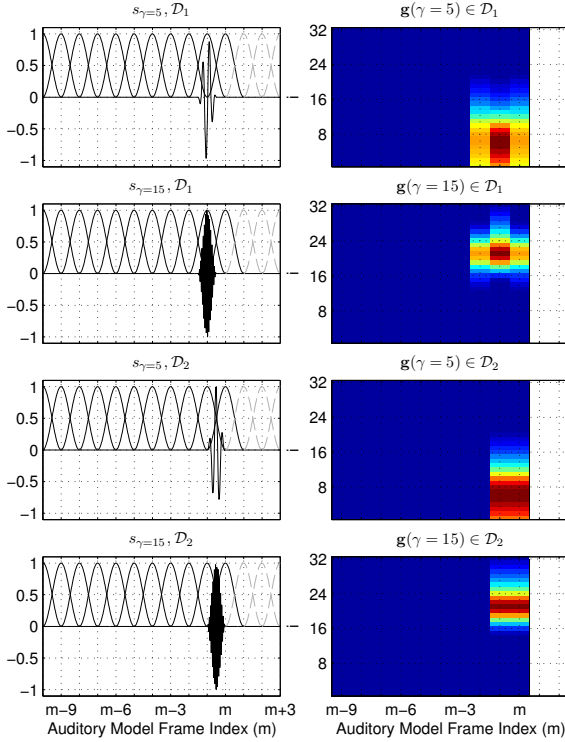
## 4. VOCODER DETAILS

CI simulations are performed with a vocoder based on sinusoidal carriers similar to [5]. In this vocoder 20 channels are used with logarithmically spaced frequencies between 150-5000 Hz. Each sinusoid is segmented into 8 ms length, 50% overlap Hann-windowed frames, which implies a channel simulation rate of 250 Hz. Note that these settings simulate the properties of the CI-processor and are chosen independently of the auditory model from Section 2.

First we will show that the time-domain additivity of the TF-spaced sinusoids in the vocoder can be preserved in the auditory domain, which validates the use of (11) in the auditory domain. Let a scaled and TF-spaced sinusoid be described as follows:

$$s_\gamma(n) = a_\gamma \cos\left(\omega_\gamma n + \phi\right) w_s(n), \quad (17)$$

where $\omega_\gamma$ and $a_\gamma$ denote the angular frequency and amplitude for channel $\gamma$, respectively, and $w_s$ its window function (the subscript $s$ of this vocoder window is used to denote its difference with the auditory model window $w_m$ from Section 2). For readability, the vocoder relevant frame-index is omitted and we assume that $w_s$ represents the current frame of interest. Since the phase is of minor importance for intelligibility in these short time frames [12], $\phi$ is assumed to be i.i.d. uniformly distributed between 0 and $2\pi$ and only the average

**Fig. 1**. Two example elements for each dictionary $\mathcal{D}_1$ and $\mathcal{D}_2$ where $\gamma = \{5, 15\}$. Left plots show realizations of $s_\gamma$ and right plots the average internal representations.
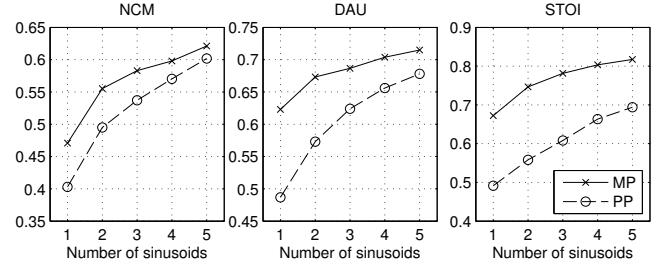
internal representation is considered. The expected value of $s_\gamma$ for one TF-unit in the auditory domain, as in (1), equals:

$$E_\phi\left[(S_\gamma)_{i,m}\right] = \frac{1}{2}\sum_k \left|\hat{h}(k)\,(\widehat{w_s e^{jn\omega_\gamma}})_m(k)\right|^2. \quad (18)$$

Moreover, the expected value of the internal representation of a sum of weighted sinusoids is given by:

$$E_\phi\left[\mathcal{I}_m\left\{\sum_\gamma a_\gamma s_\gamma\right\}\right] = \sum_\gamma |a_\gamma|^2 E_\phi\left[\mathcal{I}_m\{s_\gamma\}\right], \quad (19)$$

where the cross terms between the weighted sinusoids in the auditory domain are zero due to the i.i.d. assumption. This is a direct consequence of taking into account squared magnitudes in (1) rather than the squared root of this term. Hence, the weighted sum of sinusoids results in a squared weighted sum of average functions in the auditory domain. Note that a realization of this internal representation is expected to be close to its expected value, since the proposed metric discards all DFT-phase information in (1). Motivated by this, each element in the dictionary $\mathcal{D} = \mathbf{g}(\gamma)_{\gamma\in\Gamma}$ is defined as $\mathbf{g}(\gamma) = E[\mathcal{I}_m\{s_\gamma\}]$. The frame index $m$ is taken equal to the last frame which still overlaps with $w_s$. This means that



**Fig. 2**. Prediction results for proposed matching pursuit (MP) and peak picking (PP) algorithm (a higher score denotes more intelligible speech). The predictors STOI [3], DAU [1] and NCM [2] are all known to be reliable with vocoded speech.
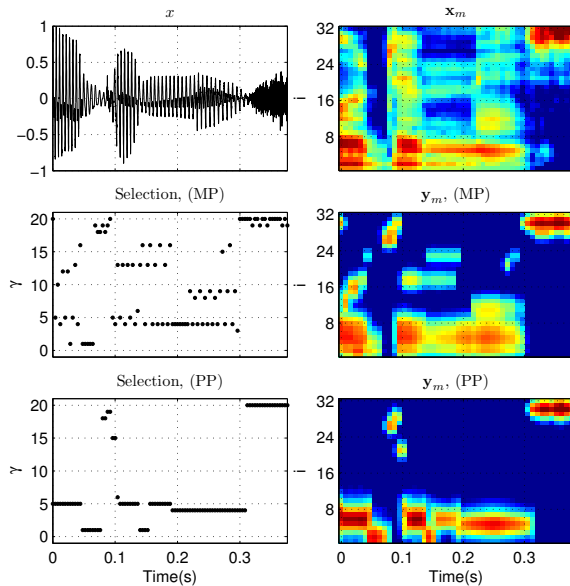
the dictionary depends on the alignment between $w_s$ and the chosen $m$. Since the support of $w_m$ (16 ms) is double the support of $w_s$ (8 ms), two possible alignments exist for which the dictionaries, say $\mathcal{D}_1$ and $\mathcal{D}_2$, can be pre-calculated and stored. Two example dictionary elements are shown for both dictionaries in Figure 1. This figure also illustrates how $m$ is chosen given $w_s$ by highlighting the windows of the auditory model. The eventual vocoded speech signal for time-frame $w_s$ is then synthesized as[1] $x \approx \sum_p \sqrt{\alpha^{(p)}} s_{\gamma^{(p)}}$.

## 5. EXPERIMENTAL RESULTS

The proposed matching pursuit (MP) algorithm is compared with the peak-picking (PP) algorithm which is currently still the basis of several existing coding strategies in CIs [4]. Signal processing details of the peak-picking algorithm can be found in [5].

Three intelligibility predictors are used to assess the intelligibility of MP and PP where the number of selected channels is varied between 1 and 5. These predictors consist of STOI [3] (the model which was simplified in Section 2), a model developed by Christiansen and Dau (DAU) [1] and the normalized covariance metric (NCM) [2]. These measures are recently proposed and can be considered as state-of-the-art for intelligibility prediction of vocoded speech. The results are shown in Fig. 2 from which we can conclude that all three measures predict that the intelligibility of MP is higher than PP. A result which is in line with informal listening tests. Largest improvements are predicted with STOI, which is not that surprising since this is the measure initially used for optimization. NCM and DAU predict that the speech intelligibility for MP with 1 sinusoid is roughly equal to the intelligibility with PP for 2 and 3 sinusoids, respectively. In the near-future real listening tests will be performed to quantify the absolute difference between MP and PP.

---

[1]In rare cases it may occur that the optimal $\alpha$ for a specific iteration is negative. Since a negative amplitude in the auditory domain does not have a meaning in the time-domain these channels are discarded.

**Fig. 3**. Auditory representations of clean and vocoded speech and channel selection for MP and PP. One channel is selected per time-instant for both algorithms.

The main differences between MP and PP are illustrated in Figure 3, where one TF-block of clean speech is used and only one channel was selected per time instant. For comparison, the clean internal representation is shown together with the internal representations for both methods, denoted by $\mathbf{y_m}$, and their corresponding channel selections. From the plots it is clear that PP tends to select the same channel independently of the previous selected channel. As a result the two formants between 0.1-0.2 seconds and channel 16-24 are completely discarded with PP, which is not the case with MP. There are two important reasons for this different behavior: (1) The proposed metric has a longer integration time such that channels selections from the past are taken into account for the current channel selection. (2) The weighting matrix $\mathbf{A}_m$ 'whitens' the speech and will therefore give a similar importance to high frequencies compared to low frequency content. Another important difference is the fact that the proposed method considers the spread over time and frequency of the sinusoids. Therefore, MP will less often select neighboring channels compared to PP.

Note that the channel stimulation rates in real CI-processors can be much higher than the rate of 250 Hz as used in the vocoder from [5]. In a real CI also the channels are typically stimulated sequentially in an interleaved manner, rather than simultaneously, in order to avoid electrical field interactions [4]. These properties of the CI cannot be included in a vocoder since no acoustical signals exist with such short-time duration and narrow frequency support. It is important to add, however, that these are constraints of the use of any vocoder and not of the proposed channel selection method.

Namely, the dictionary can be easily extended to shorter pulse durations in a real CI environment.

## 6. CONCLUDING REMARKS

In this paper it is shown that the existing short-time objective intelligibility (STOI) measure can be expressed as a weighted $\ell_2$ norm in the auditory domain. Due to the mathematical properties of this norm it facilitated the use of the matching pursuit algorithm in the channel selection technique in cochlear implants (CIs). Acoustic CI simulations are generated based on a sinusoidal vocoder where a large intelligibility improvement was found by three state-of-the-art intelligibility predictors compared to a peak-picking algorithm.

## 7. REFERENCES

[1] C. Christiansen, M. S. Pedersen, and T. Dau, "Prediction of speech intelligibility based on an auditory preprocessing model," *Speech Communication*, vol. 52, pp. 678–692, 2010.

[2] F. Chen and P. Loizou, "Predicting the intelligibility of vocoded speech," *Ear and Hearing*, vol. 32, no. 3, p. 331, 2011.

[3] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, "An algorithm for intelligibility prediction of time-frequency weighted noisy speech," *IEEE Trans. Audio Speech Lang. Process.*, vol. 19, no. 7, pp. 2125–2136, 2011.

[4] B. Wilson and M. Dorman, "Cochlear implants: a remarkable past and a brilliant future," *Hearing research*, vol. 242, no. 1-2, pp. 3–21, 2008.

[5] M. F. Dorman, P. C. Loizou, A. J. Spahr, and E. Maloff, "A comparison of the speech understanding provided by acoustic models of fixed-channel and channel-picking signal processors for cochlear implants," *J Speech Lang Hear Res*, vol. 45, no. 4, pp. 783–788, 2002.

[6] P. Seligman, H. McDermott *et al.*, "Architecture of the spectra 22 speech processor," *Annals of Otology, Rhinology and Laryngology*, vol. 104, no. suppl 166, pp. 139–141, 1995.

[7] W. Nogueira, A. Büchner, T. Lenarz, and B. Edler, "A psychoacoustic nofm-type speech coding strategy for cochlear implants," *EURASIP Journal on Applied Signal Processing*, vol. 2005, pp. 3044–3059, 2005.

[8] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, "An evaluation of objective measures for intelligibility prediction of time-frequency weighted noisy speech," *J. Acoust. Soc. Am.*, vol. 130, no. 5, pp. 3013–3027, 2011.

[9] R. Drullman, J. Festen, and R. Plomp, "Effect of temporal envelope smearing on speech reception," *J. Acoust. Soc. Am.*, vol. 95, no. 2, pp. 1053–1064, 1994.

[10] S. Mallat and Z. Zhang, "Matching pursuits with time-frequency dictionaries," *IEEE Transactions on Signal Processing*, vol. 41, no. 12, pp. 3397–3415, 1993.

[11] S. van de Par, A. Kohlrausch, R. Heusdens, J. Jensen, and S. Jensen, "A perceptual model for sinusoidal audio coding based on spectral integration," *EURASIP J. on Appl. Signal Processing*, vol. 2005, no. 9, pp. 1292–1304, 2005.

[12] L. Liu, J. He, and G. Palm, "Effects of phase on the perception of intervocalic stop consonants," *Speech Communication*, vol. 22, no. 4, pp. 403 – 417, 1997.