# A NEW FRAMEWORK FOR CONVERGENCE ANALYSIS AND ALGORITHM DEVELOPMENT OF ADAPTIVE IIR FILTERS

*Phillip M. S. Burt* *

Dept. PTC
Escola Politécnica
Universidade de São Paulo
CEP 05508-900 São Paulo SP Brazil

*Phillip A. Regalia*

Dept. CITI
CNRS/SAMOVAR UMR 5157
Institut National des Télécommunications/GET
9 rue Charles Fourier
91011 Evry cedex France

## ABSTRACT

A parameterization of an adaptive IIR filter's poles is developed, based on balanced realization theory. From this we develop a local approximation of the actual adapted pole parameters, in which convergence speed is related to a certain eigenvalue spread. This, in turn, is shown to relate to the Hankel singular values of the system to be identified, as well as certain coefficient sensitivity functions of the adapted filter. Based on these properties, a new adaptive IIR algorithm is proposed. In order to achieve faster convergence, it combines an adaptive lattice with function approximation.

## 1. INTRODUCTION

Adaptive IIR filters can in principle represent an advantageous alternative in relation to adaptive FIR filters, due to their capacity of providing long impulse responses with a small number of coefficients. Problems related to local minima, stability and the effect of poles close to the unit circle have been tackled, over the years, by several authors, leading to different adaptive algorithms and realization structures. In most cases, the numerical examples that are given are second order cases.

One aspect of adaptive IIR filters that hasn't received much attention is that the perfomance of simple constant gain algorithms can rapidly degrade as the order of the filter grows. It can easily be verified by simulations [1] that, even in the absence of local minima, the adaptation of filters with order greater than 2 (4, for instance) can remain almost stopped in regions were the mean square error is far from being acceptable. While other adaptive algorithms, such as the Newton kind, are less susceptible to this effect, their greater computational complexity undermines the main motivation of adaptive IIR filters, which is low computational complexity. More insight into this aspect of the

convergence of adaptive IIR filters is, therefore, desirable. This is the problem we address here.

## 2. ANALYTICAL FRAMEWORK

### 2.1. True gradient adaptation on the reduced error surface

We consider, initially, that a rational function $\widehat{H}(z)$ is adapted so as to minimize the mean square error between the output $\widehat{y}(n) = \widehat{H}(z)u(n)$ produced for a known white input $u(n)$ and the output of a system $H(z)$ to the same input, $y(n) = H(z)u(n) + \eta(n)$. In this mixed notation, $z$ is the unit-delay operator, with $zu(n) = u(n-1)$. Assuming additive noise $\eta(n)$ is independent of $u(n)$ makes the problem equivalent to the minimization of the norm $||H(z) - \widehat{H}(z)||^2$.

A less general problem, now, is to assume $H(z)$ is known. If the order imposed on $\widehat{H}(z)$ is smaller than the order of $H(z)$, this problem doesn't have a closed-form solution and is still an optimization problem of interest. One possible procedure for solving it would be as follows. 1) Choose an initial value for the poles of $\widehat{H}(z)$. 2) Given these poles, obtain the zeros of $\widehat{H}(z)$ that minimize $||H(z) - \widehat{H}(z)||^2$(this problem has a closed-form solution). The error thus achieved, denoted by $||g(z)||^2 = ||H(z) - \widehat{H}(z)||^2$, depends only on the poles of $\widehat{H}(z)$ and is termed the *reduced error surface*. 3) Adapt whatever are the pole parameters $w_k$ by the gradient, so that at iteration $n+1$ we have

$$w_k(n+1) = w_k(n) - \frac{\mu}{2}\frac{\partial}{\partial w_k}\|g(z)\|^2, \qquad (1)$$

and go back to step 2. The fact that $H(z)$ is known here by no means makes the aforementioned problem of very slow convergence disappear. For our purposes this is good: as will be seen, when (1) is cast in a different parameterization a powerful insight is gained in relation to the origin of this convergence problem. Moreover, this extends to our original problem, where

$H(z)$ is not known and stochastic gradient adaptation is employed, and where the assumption that the orders of $\widehat{H}(z)$ and $H(z)$ are equal doesn't make the problem trivial.

## 2.2. SVD analysis of the reduced error surface

We consider that in $\widehat{H}(z)$ we can vary $M$ zeros, $M$ poles and a gain. We say, therefore, that $\widehat{H}(z)$ is "of order $M$". It can be shown that, in this case, on the reduce error surface we have $H(z) - \widehat{H}(z) = g(z)V(z)$, where $g(z)$ is strictly causal and $V(z)$ is the unit-norm all-pass function that has the same poles as $\widehat{H}(z)$ [1]. From this it also follows that

$$||g(z)||^2 = \mathbf{v}^t \Gamma_H^2 \mathbf{v}, \qquad (2)$$

where $\Gamma_H$ is the Hankel form of system $H(z)$ and vector $\mathbf{v}$ contains the coefficients of the expansion of $V(z)$. It should be noted that $||g(z)||^2 = ||H(z) - \widehat{H}(z)||^2$, since $V(z)$ is all-pass.

Now, if $H(z)$ has order $N$ then $\Gamma_H$ can be written as a function of its singular values $\sigma_k$ and Schmidt pairs $(\boldsymbol{\zeta}_k, \boldsymbol{\eta}_k)$ as $\Gamma_H = \sum_{k=1}^{N} \boldsymbol{\zeta}_k \sigma_k \boldsymbol{\eta}_k^t$ [2], where sets $\boldsymbol{\eta}_k$ and $\boldsymbol{\zeta}_k$ are orthonormal. From this and (2) it results that $||g(z)||^2 = \sum_{k=1}^{N} \mathbf{v}^t \boldsymbol{\zeta}_k \sigma_k^2 \boldsymbol{\zeta}_k^t \mathbf{v}$. We can then return to polynomial form and write

$$||g(z)||^2 = \sum_{k=1}^{N} \sigma_k^2 \langle \zeta_k(z), V(z) \rangle^2 = \sum_{k=1}^{N} \sigma_k^2 \alpha_k^2, \qquad (3)$$

where $\alpha_k \doteq \langle \zeta_k(z), V(z) \rangle \leq 1$. Functions $\zeta_k(z)$ are scaled (to unit norm) controllability functions of a balanced realization of $H(z)$. As can be seen, the error is a quadratic function of the terms $\alpha_k$. In the following, we consider $N = M$ and show that the terms $\alpha_k$ can also be used to parameterize the poles of $\widehat{H}(z)$.

## 2.3. Balanced form linked parameterization

We consider initially a function $\tilde{H}(z) = C(z)/D(z)$, of order $M$ as $\widehat{H}(z)$. $\tilde{H}(z)$ is not necessarily equal to the system $H(z)$ considered above. In the following, also, the notation $\overline{F}(z) \doteq z^M F(z^{-1})$ is used, for any $F(z)$. It can be shown that the $M$ scaled controllability functions $\tilde{\zeta}_k(z)$ of a balanced implementation of $\tilde{H}(z)$ and the functions $\{U(z), zU(z), z^2 U(z), \ldots\}$, $U(z) \doteq \overline{D}(z)/D(z)$, constitute an orthonormal basis for the space $H_2$ of causal and stable functions of $z$ [1]. From the denominator of $\widehat{H}(z) = B(z)/A(z)$ we form the all-pass function $V(z) = \overline{A}(z)/A(z)$, which belongs to $H_2$. Therefore, with $P(z)/Q(z) \in H_2$, $V(z)$ can always be written in function of that basis as

$$V(z) = \frac{\overline{A}(z)}{A(z)} = \sum_{k=1}^{M} \alpha_k \tilde{\zeta}_k(z) + \frac{P(z)}{Q(z)} \frac{\overline{D}(z)}{D(z)}. \qquad (4)$$

If the $\alpha_k$ are given, then it can be shown that (4) leads to a system of linear equations, where the $M$ coefficients of monic $A(z)$ and the $M + 1$ coefficients of $P(z)$ make up a total of $2M + 1$ unknowns, for the same number of equations. We have, then,

$$\left[\begin{array}{cc} \theta_D - \theta_R T & -\theta_{\overline{D}} \end{array}\right] \left[\begin{array}{c} \mathbf{a} \\ \mathbf{p} \end{array}\right] = \mathbf{b},$$

where, for any $F(z)$, $\theta_F$ is a convolution matrix composed of the coefficients of $F(z)$, $T$ is the anti-diagonal permutation matrix, and $\mathbf{a}$, $\mathbf{p}$ and $\mathbf{b}$ contain the coefficients of $A(z)$, $P(z)$ and $R(z) - z^M D(z)$, respectively.

If, otherwise, the coefficients of $A(z)$ are given, the parameters $\alpha_k$ can be obtained directly from $\alpha_k = \langle \tilde{\zeta}_k(z), V(z) \rangle$. Alternatively, (4) gives a system of $2M + 1$ linear equations and the same number of unknowns,

$$\left[\begin{array}{cc} \theta_A \Phi_R & \theta_{\overline{D}} \end{array}\right] \left[\begin{array}{c} \boldsymbol{\alpha} \\ \mathbf{p} \end{array}\right] = \mathbf{c},$$

where $\boldsymbol{\alpha}$ contains the parameters $\alpha_k$, $\mathbf{c}$ has the coefficients of $\overline{A}(z)D(z)$ and $\Phi_R$ is composed of the coefficients of the numerators of $\tilde{\zeta}_k(z)$.

As seen, then, we can go back and forth between the coefficients of $A(z)$ and parameters $\alpha_k$, which constitute, therefore, a different parameterization of the poles of $\widehat{H}(z) = B(z)/A(z)$. This parameterization is linked to the function $\tilde{H}(z)$ adopted in the beginning.

It may be possible that an adaptive method can be developed where parameters $\alpha_k$ are, themselves, the adapted parameters. In this scheme, the linked function $\tilde{H}(z)$ would play the role of an initial estimate of $H(z)$. The motivation for this would be the fact that if $\tilde{H}(z) \approx H(z)$ then the error surface would be close to quadratic and, as follows from (3), the problem of slow convergence could be overcome by simply using a different adaptation gain $\mu_k = \mu/\sigma_k^2$ for each parameter. This idea is not pursued here, however. Otherwise, in the following we will use parameters $\alpha_k$ to *describe* the adaptation process of a given set of parameters $w_k$ by the gradient. Besides the greater understanding of the convergence of gradient algorithms this wil provide, it will also lead to an algorithm that attempts to overcome the problem of slow convergence.

## 2.4. Local approximation

We are now in a position to return to the adaptation given by (1). As follows from the previous discussion, the poles of $\widehat{H}(z)$ can be parameterized by the balanced form linked parameters $\alpha_k$, which we group in $\boldsymbol{\alpha} = [\alpha_1 \ldots \alpha_M]^t$. Moreover, when the linked function satisfies $\tilde{H}(z) = H(z)$, the error is entirely determined by $\boldsymbol{\alpha}$ as given by (3). Therefore, the adaptation of the parameters $\mathbf{w} = [w_1 \ldots w_M]^t$ used for the poles, and the resulting error, are entirely described by the

sequence $\boldsymbol{\alpha}(n)$. Of course, for this to be useful it is necessary in the first place to be able to write $\boldsymbol{\alpha}(n+1)$ in function of $\boldsymbol{\alpha}(n)$. To this end, a first-order approximation of $\boldsymbol{\alpha}(n+1)$ is given by

$$\tilde{\boldsymbol{\alpha}}(n+1) = \boldsymbol{\alpha}(n) + \left.\frac{\partial\boldsymbol{\alpha}}{\partial\mathbf{w}}\right|_{\mathbf{w}(n)} [\mathbf{w}(n+1) - \mathbf{w}(n)], \quad (5)$$

where $\left.\frac{\partial\boldsymbol{\alpha}}{\partial\mathbf{w}}\right|_{\mathbf{w}(n)} \doteq \mathbf{J}(\boldsymbol{\alpha}(n))$ is the Jacobian matrix of derivatives ( $\partial\alpha_i/\partial w_j$ at row $i$ and column $j$) at point $\mathbf{w}(n)$. Now, considering adaptation of $\mathbf{w}$ by the gradient, from (1) we have $\mathbf{w}(n+1) - \mathbf{w}(n) = -\frac{\mu}{2}\left.\frac{\partial\varepsilon}{\partial\mathbf{w}}\right|_{\mathbf{w}(n)}$, where we use $\varepsilon \doteq ||g(z)||^2$, for notational simplicity. With (5) this leads to

$$\tilde{\boldsymbol{\alpha}}(n+1) = \boldsymbol{\alpha}(n) - \frac{\mu}{2}\left.\frac{\partial\boldsymbol{\alpha}}{\partial\mathbf{w}}\frac{\partial\varepsilon}{\partial\mathbf{w}}\right|_{\mathbf{w}(n)} .$$

Applying the chain law of differentiation, we can always write $\frac{\partial\varepsilon}{\partial\mathbf{w}} = \left(\frac{\partial\boldsymbol{\alpha}}{\partial\mathbf{w}}\right)^t \frac{\partial\varepsilon}{\partial\boldsymbol{\alpha}}$. Also, from (3), we have $\frac{\partial\varepsilon}{\partial\boldsymbol{\alpha}} = 2\boldsymbol{\Sigma}^2\boldsymbol{\alpha}$, where $\boldsymbol{\Sigma}$ is the diagonal matrix with singular values $\sigma_k$. We arrive, finally, at

$$\tilde{\boldsymbol{\alpha}}(n+1) = \left[\mathbf{I} - \mu\mathbf{J}(\boldsymbol{\alpha}(n))\mathbf{J}(\boldsymbol{\alpha}(n))^t\boldsymbol{\Sigma}^2\right]\boldsymbol{\alpha}(n). \quad (6)$$

This expression applies at any point $\boldsymbol{\alpha}(n)$, no assumption that $\boldsymbol{\alpha}(n)$ is a stationary point of the adaptation having been made. As can be seen, there is a partial separation of the effect of the system to be identified, which affects $\boldsymbol{\Sigma}$ and $\mathbf{J}$, and the effect of the parameterization of the adaptation, which affects only $\mathbf{J}$. The expression is also similar to what is obtained for the steepest-descent algorithm in adaptive FIR filtering, the eigenvalues of $\mathbf{J}(\boldsymbol{\alpha}(n))\mathbf{J}(\boldsymbol{\alpha}(n))^t\boldsymbol{\Sigma}^2$ playing the part of the eigenvalues of the input correlation matrix. These eigenvalues are non-negative, which can be seen noting that the non-null eigenvalues of $\mathbf{JJ}^t\boldsymbol{\Sigma}^2$ and $\mathbf{J}^t\boldsymbol{\Sigma}^2\mathbf{J} = (\boldsymbol{\Sigma}\mathbf{J})^t(\boldsymbol{\Sigma}\mathbf{J})$ are always equal. The larger of them will tend to limit the value of the gain $\mu$. Convergence will be slow then at a point $\boldsymbol{\alpha}(n)$ if the minimum eigenvalue of $\mathbf{J}(\boldsymbol{\alpha}(n))\mathbf{J}(\boldsymbol{\alpha}(n))^t\boldsymbol{\Sigma}^2$ is small and $\boldsymbol{\alpha}(n)$ is in the direction of the associated eigenvector.

### 2.5. Relation with Hessian approximation

Usually, convergence analysis of adaptive IIR filters is carried out directly in terms of the adapted parameters and local approximation is restricted to stationary points $\mathbf{w}_*$. Writing the Hessian matrix as a function of the parameters $\alpha_k$, it can be shown for gradient adaptation that at a stationary point $\mathbf{w}_*$ we have $\mathbf{w}(n+1) - \mathbf{w}_* \approx \left[\mathbf{I} - \mu\mathbf{J}^t(\mathbf{w}_*)\boldsymbol{\Sigma}^2\mathbf{J}(\mathbf{w}_*)\right][\mathbf{w}(n) - \mathbf{w}_*]$. As in the local approximation (6), the use of the balanced form linked parameterization leads to a partial separation of the effects of the system to be identified and the choice of the adapted parameters. The local approximation, though, is more useful since it is not restricted to stationary points.
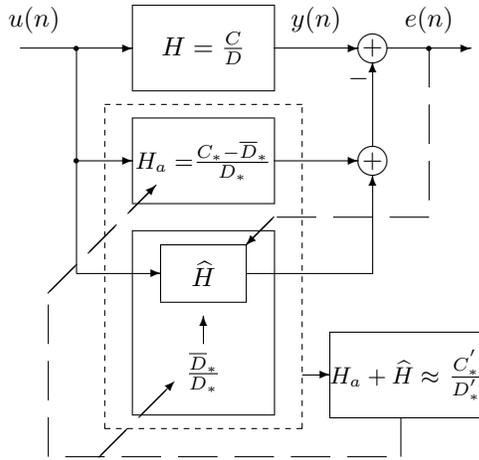
### 3. JACOBIAN GRAMMIAN PROPERTIES

For analyzing the eigenvalue spread of $\mathbf{JJ}^t\boldsymbol{\Sigma}^2$, denoted $\chi_{J\Sigma}$, it can be verified that the product upper bound $\chi_J\chi_\Sigma$ is often useful. $\chi_\Sigma$ depends on the relative position of the poles and zeros of $H(z)$, being equal to 1 when $H(z)$ is all-pass. To deal with $\chi_J$, we call upon the orthonormal controllability functions $\xi_k(z) = \hat{D}_k(z)/D(z)$ of a lattice realization of $H(z)$ and define $\beta_k \doteq \langle\xi_k(z), V(z)\rangle$. We have $\mathbf{J}_\beta = \mathbf{QJ}$ for the associated Jacobian, with $\mathbf{Q}$ orthogonal, and, therefore, $\chi_J = \chi_{J,\beta}$. This shows that $\chi_J$ depends only on the poles of $H(z)$. When the adapted parameters are the direct form parameters, it can be also be shown that $[\mathbf{J}_\beta^{-t}\mathbf{J}_\beta^{-1}]_{i,j} = \langle\hat{D}_i(z), \hat{D}_j(z)\rangle$ at the global minimum. Using subscript $o$ in association with the global minimum, it follows that $\chi_{J,o}$ is closer to one when the poles of $H(z)$ are away from the unit circle and/or uniformly distributed in angle. A trivial case is $D(z) = 1 - a^M z^M$, which gives $\chi_{J,o} = 1$. For the lattice form only partial analytical properties have been obtained so far. However, it has been observed numerically that, in relation to the direct form case, $\chi_{J,o}$ tends to be closer to one.

### 4. SUCESSIVE APPROXIMATIONS

An adaptive IIR algorithm that makes use of the properties of matrices $\mathbf{JJ}^t$ and $\boldsymbol{\Sigma}$ seen in the previous section is presented in the following. Due to space limitation, the algorithm is not described in a more formal manner and some details are left out. A block diagram of the algorithm is at the end of this paragraph. From an initial estimate $C_*(z)/D_*(z)$ of system $H(z) = C(z)/D(z)$ (estimate that in some cases may be the origin of the coefficient space), the transfer function of a fixed auxiliary block is set at $H_a(z) = [C_*(z) - \overline{D}_*(z)]/D_*(z)$ and an $M$ order lattice adaptive filter is initialized with $\widehat{H}(z) = \overline{D}_*(z)/D_*(z)$. The simplified partial (stochastic) gradient algorithm [1] is employed. Assuming $H(z)$ and $H_a(z)$ also have order $M$, the overall system $H(z) - H_a(z)$ to be identified by the lattice has order $2M$. Depending on the initial estimate, the overall system will be close to all-pass and the first $M$ eigenvalues of $\boldsymbol{\Sigma}$ will be close to one. We assume that in this case the adaptive lattice $\widehat{H}(z)$ will converge rapidly, say in $n_a$ iterations, to a point were the first $M$ error terms in (3) are close to zero. (This assumption is based on the analysis presented in the previous sections, though the analysis has yet to be extended to the undermodelled case.) After $n_a$ iterations, then, the function now given by $\widehat{H}(z) + H_a(z)$, which is of order $2M$, is approximated by a new $M$ order function $C_*(z)/D_*(z)$ and the whole process starts again. Different closed-form approximation methods can be employed for the last step, which is directly related to model reduction problems (e.g.,

[3]). Here, we have used minimization of the equation error $||H(z)D_*(z) - C_*(z)||^2$. An important point is that the approximation procedure is not performed at each iteration. In particular, we consider that after the $n_a$ adaptation iterations, $n_x$ iterations are used to compute the approximation, period during which the adapted parameters are frozen. If $n_x$ is large enough, the per-iteration computational load of the approximation will be small. Another important point is that after introducing the new approximation the adapted parameters remain frozen still for $n_w$ iterations, so that the error transient can decay.
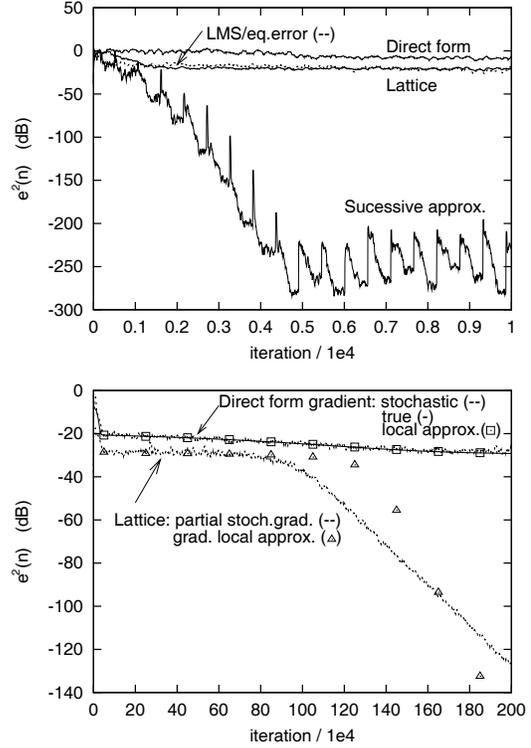


## 5. NUMERICAL RESULTS

One $H(z^{-1})$ was selected from of a set of 100 randomly generated fourth-order unit-norm transfer functions. Its zeros are at $0.20 \pm 1.60j$ and $-0.051 \pm 0.81j$, and its poles are at $0.47 \pm 0.52j$ and $0.78 \pm 0.081j$. The sucessive approximations (SA) algorithm, with $n_a = 300$, $n_x = 200$ and $n_w = 50$ was compared with the more conventional direct form (stochastic gradient) and lattice (simplified partial stochastic gradient) cases. As equation error approximation was used in the SA algorithm, an LMS/equation error algorithm was also considered for comparison. For each algorithm the employed step $\mu$ was within a factor of 2 of its maximum value. Adaptation started from the origin in all cases. No output noise was considered in order to highlight the previously analyzed convergence properties. The results for one realization of the white input are in the first figure that follows. Smoothing was employed for better visualization. It can be seen that the proposed algorithm converges much faster than the other cases.

Results of a longer simulation are in the second figure. For the direct form, the stochastic gradient algorithm is closely approximated by the true gradient algorithm given by (1) and the local approximation given by $\tilde{\alpha}(n+L) = \left[ \mathbf{I} - \mu \mathbf{J}(\tilde{\alpha}(n))\mathbf{J}(\tilde{\alpha}(n))^t \mathbf{\Sigma}^2 \right]^L \tilde{\alpha}(n)$, with $L = 10^4$ ($\tilde{\alpha}(n)$ is only shown every $2 \times 10^5$ iterations).

For the lattice, a true gradient algorithm was not implemented, due to its complexity. The local approximation, with $L = 10^3$, initially follows closely but then deviates from the partial gradient algorithm. This is not unexpected, as the latter employs a biased approximation of the gradient in order to reduce computational complexity.





For the considered $H(z)$, the eigenvalue spread of $\mathbf{\Sigma}^2$ is 36 dB. The other eigenvalue spreads (in dB) and the product bound are in the table below, calculated at the global minimum $\widehat{H}(z) = H(z)$ and at the points reached after $5 \times 10^4$ iterations. It can be seen that at these points the eigenvalue spread is considerably higher than at the global minimum, which exemplifies the greater utility of the local approximation in comparison with the Hessian approximation.

| Form | $\widehat{H} = H$ | | | $n = 5 \times 10^4$ | | |
|---|---|---|---|---|---|---|
| | $\chi_J$ | $\chi_J \chi_\Sigma$ | $\chi_{J\Sigma}$ | $\chi_J$ | $\chi_J \chi_\Sigma$ | $\chi_{J\Sigma}$ |
| Direct | 36 | 72 | 67 | 49 | 85 | 79 |
| Lattice | 7 | 43 | 36 | 29 | 64 | 62 |

## 6. REFERENCES

[1] P. A. Regalia, *Adaptive IIR filtering in signal processing and control*, Marcel Dekker, New York, 1995.

[2] F. Deprettere, Ed., *SVD and signal processing*, North-Holland, Amsterdam, 1988.

[3] R. Johansson, "System identification using LQG-balanced model reduction," in *Proc. 41st IEEE Conf. on Decision and Control*, 2002, pp. 258–263.