

Optimal Training for Residual Self-Interference for Full Duplex One-way Relays

Xiaofeng Li, Cihan Tepedelenlioğlu *Member, IEEE*,
and Habib Şenol *Member, IEEE*

Index Terms

Full duplex relays, residual self-interference, maximum likelihood estimation, optimal training sequence, Toeplitz matrix, frequency selective channels, multiple relays

Abstract

Optimal training design and maximum likelihood channel estimation for one-way amplify-and-forward full duplex relay systems are proposed. The destination estimates the residual self-interference (RSI) channel as well as the end-to-end channel of the relay system, aiming to cancel the RSI through equalization. The log-likelihood function is maximized through a quasi-Newton method, with an MMSE-based initialization. The Cramer-Rao bounds are derived to evaluate the accuracy of the estimates. By using Szegő's theorems about Toeplitz matrices, we minimize the Cramer-Rao bounds and find the corresponding optimal training sequences. Extensions of our method to frequency selective channels and multiple relays are also presented.

I. INTRODUCTION

Due to the growing demands on wireless bandwidth, the need for high spectral efficiency has become more urgent. In-band full duplex (FD) relays, which transmit and receive simultaneously in the same frequency band, offer a viable solution since they theoretically have the ability to double the spectral efficiency, compared to half duplex relays. In practice, however, the relay receives

Xiaofeng Li and Cihan Tepedelenlioğlu are with the Department of Electrical, Computer and Energy Engineering, Arizona State University, Tempe, AZ, 85287, USA (email: xiaofen2, cihan@asu.edu).

Habib Şenol is with Department of Computer Engineering, Faculty of Engineering and Natural Sciences, Kadir Has University, Istanbul 34083, Turkey (e-mail: hsenol@khas.edu.tr).

strong self-interference in FD mode, which is challenging to overcome. Recently, self-interference cancellation techniques have been developed with great promise [1]–[3]. With these techniques, the self-interference can be canceled by estimating the self-interference channels [4]–[6] in FD mode, or be suppressed with null-space methods in MIMO systems [7]. However, despite these advances, residual self-interference (RSI) still exists after the self-interference cancellation [3], [7]–[9]. Therefore accurate channel estimation in the presence of RSI is required at the destination to further improve the system performance by canceling RSI.

Several works analyze the system performance in the presence of RSI with different criteria such as interference power, outage probability, and bit error rate (BER) [9]–[12]. Reference [10] investigates the outage probability of an amplify-and-forward (AF) FD relay. Reference [9] takes advantages of multiple antennas to suppress the RSI power using a null-space pre-coding matrix. References [11] and [12] analyze the diversity and the capacity of FD relays in the presence of RSI. However, these works do not consider the cancellation of the RSI. Some of the works assume perfect channel state information (CSI) [10], [11] while others assume imperfect CSI [9] without considering how to estimate the channels. As shown by the results of these works, the system performance suffers from the RSI since it is still quite high as compared to the desired received signal, and does not yield good performance when treated as noise. In [10], the authors assume that the RSI power is as high as 10 dB than the desired signal power and also results in inter-symbol interference (ISI) in AF relays. References [3], [8], and [9] report that the power of RSI may be as high as 30 dB than the noise floor even after applying self-interference cancellation. This motivates incorporating the RSI into the end-to-end channel model, and then estimating it, and removing it at the receiver.

In addition to the one-way relay system, in-band FD can also be applied in two-way relays in which two sources exchange their messages through the relay [13]–[17]. The achievable rates of FD relays and half duplex relays are compared in [13]. References [14], [15] analyze the effect of channel estimation error and suppression of the self-interference without showing how to estimate the channels. Reference [16] designs precoder matrices to suppress the self-interference and analyzes the end-to-end SNR performance in the absence of RSI. The channel estimation problem for RSI

in FD two-way relays is addressed in [17]. The RSI channel and the channels between nodes are estimated simultaneously at the sources, and the asymptotic behavior of the Fisher information is analyzed.

From the RSI channel estimation point of view, reference [18] proposes two methods for the RSI channel in an AF relay system where the destination is equipped with massive MIMO. In the first method the authors consider the case where the RSI channel is estimated by the relay itself, and in the second method the base station estimates the RSI channel. However in their system model, the RSI is incorporated with the noise term and ISI caused by the RSI is treated as noise. A conference version of this manuscript [19] proposes a maximum likelihood estimator for the RSI channel and derives the Cramer-Rao bounds (CRBs). In this work, we further analyze the CRB by using asymptotic properties of Toeplitz matrices and design the optimal training sequences. We also extend our training method to the case in which the channels in the relay system are frequency selective and the case of multi-relay systems, which were not considered in [19]. Reference [17] investigates the channel estimation problem in FD two-way relay which is fundamentally different than the one-way relay setting considered herein. In addition to the difference in the relaying technique (one-way versus two-way) which yields a fundamentally different system model, this work analyzes the CRB which is one step further than the analysis of Fisher information in [17] because of the tractability of the inverse of the latter. Furthermore, the optimization of the training sequence is considered, exploiting the closed-form expression of the CRB. Moreover, we also consider important generalizations such as the multiple relay scenario and the frequency-selective case which are not considered in [17].

In this work, we consider an AF FD one-way relay system with the relay working in FD mode. The RSI in the relay propagates to the destination, creating an end-to-end ISI channel. We further cancel the RSI at the destination by estimating the RSI channel and applying equalization. Different from studies which focus on canceling the self-interference at the relay itself, we reduce the complexity burden at the relay and aim to cancel the interference at the destination. Thus, the destination performs the estimation and cancellation operations. To estimate the RSI channel as well as the channels between nodes, a maximum likelihood (ML) estimator is proposed, which

is formulated by maximizing the log-likelihood function through a quasi-Newton method. The quasi-Newton method is initialized by a linear minimum mean square error (MMSE) estimator. The CRBs are derived to evaluate the accuracy of the estimates. By using Szegő's theorems for asymptotic Toeplitz matrices, we are able to find the corresponding optimal training sequences through minimizing the CRB. We also extend our training method to the case when the channels between nodes are frequency selective and the case of multiple relay systems.

The rest of the paper is organized as follows: Section II describes the system model of an AF FD one-way relay system. Section III proposes the training scheme and the ML estimator. The CRBs are derived and analyzed in Section IV, and the optimal and the approximately optimal training sequences are also discussed. Section V and VI extend our training method to frequency selective channel and multi-relay case. Section VII shows the numerical results and Section VIII concludes the paper.

II. SYSTEM MODEL

We consider a system consisting of a source, a relay, and a destination, without any direct link between the source and the destination. Amplify-and-forward relay protocol is adopted. The relay uses two antennas, one receiving the current symbol while the other one amplifying and forwarding the previously received symbol, to operate in FD mode [10]. The channel coefficients between the source and the relay, and the relay and the destination are h_{sr} and h_{rd} respectively. The two channels between nodes are assumed to be flat fading modeled by independent complex Gaussian random variable with zero means and variances σ_{sr}^2 , σ_{rd}^2 , respectively. A separate pre-stage is assumed to gather the information of the self-interference channel to perform analog and digital cancellation methods in the next transmission stage [3], [18]. During the transmission, the self-interference is reduced in RF with analog cancellation methods until the residual self-interference power falls in the ADC dynamic range, and then is further suppressed by digital methods. However, despite these suppression methods, the RSI is still present at the destination. Since the line of sight (LoS) component of the self-interference varies very slowly, it can be much reduced by the RF cancellation, which means the reflected multi-path component dominates the RSI in the transmission stage. Therefore, in our system the RSI channel h_{rr} is modeled as a complex Gaussian random variable

with zero mean and variance σ_{rr}^2 . This is different from reference [8] because it considers the self-interference channel with the LoS component between the transmitting and receiving antennas (not the residual) and models the channel as Rician fading with a low K-factor. The Gaussian RSI model mentioned above is also used in other works such as [16], [18], [20].

Even though the LoS component is largely canceled, the RSI power is still not small enough to be treated as noise, and often higher than the desired signal power. Moreover, the RSI makes the overall end-to-end channel an ISI channel in the AF relay even when the channels on all links are flat fading. The CSI of the RSI channel is needed for equalizers to alleviate the ISI at the destination receiver.

Following [10], we assume there is one symbol processing delay for the relay to forward its received symbols. Let the transmitted signal at the relay be $t_{\text{r}}[n] = \alpha y_{\text{r}}[n - 1]$. At the destination, the received symbol at the n th time interval is

$$\begin{aligned} y_{\text{d}}[n] &= h_{\text{rd}}t[n] + n_{\text{d}}[n] = h_{\text{rd}}(\alpha y_{\text{r}}[n - 1]) + n_{\text{d}}[n] \\ &= \sum_{k=1}^{\infty} h\theta^{k-1}x[n - k] + \sum_{k=1}^{\infty} d\theta^{k-1}n_{\text{r}}[n - k] + n_{\text{d}}[n] \quad n = 0, 1, \dots \end{aligned} \quad (1)$$

where $y_{\text{r}}[n] = h_{\text{sr}}x[n] + \alpha h_{\text{rr}}y_{\text{r}}[n - 1] + n_{\text{r}}[n]$ is the n th received symbol at the relay and $x[n]$ is the transmitted signal of the source. For brevity, we define $d := \alpha h_{\text{rd}}$, $h := \alpha h_{\text{sr}}h_{\text{rd}}$ and $\theta := \alpha h_{\text{rr}}$. Noise terms $n_{\text{r}}[n]$ and $n_{\text{d}}[n]$ are complex Gaussian with zero mean and variance σ_{n}^2 . If there is no RSI, the effective end-to-end channel h is the overall channel for the system. However, the self-interference link θ forms a feedback at the relay, which makes the overall channel a single pole infinite impulse response (IIR) channel and causes ISI. Additionally, the effective noise at the destination is colored with correlations that depend on the pole.

The self-interference cancellation at the relay should be such that $|\theta| < 1$ is possible with proper choice of α . Such α keeps the system stable and guarantees finite average relay transmit power. The average relay transmit power is calculated as

$$\mathbb{E}[t_{\text{r}}[n]t_{\text{r}}^*[n]] = \alpha^2 \sum_{k=1}^{\infty} (\alpha^2 |h_{\text{rr}}|^2)^{(k-1)} (P_{\text{s}}|h|^2 + \sigma_{\text{n}}^2) = \alpha^2 \frac{P_{\text{s}}|h|^2 + \sigma_{\text{n}}^2}{1 - \alpha^2 |h_{\text{rr}}|^2}. \quad (2)$$

The condition for the stability of the system is given by [21]

$$\mathbb{E}[t_r[n]t_r^*[n]] \leq P_r, \quad (3)$$

where the expectation in (3) is with respect to the noise. By solving (3), α should satisfy $\alpha^2|h_{rr}|^2 = |\theta|^2 < 1$. However, during channel estimation, the instantaneous CSI is not known to the relay. We can choose α to satisfy a long term condition $\mathbb{E}[\alpha^2|h_{rr}|^2] < 1$ which leads to $\alpha^2\sigma_{rr}^2 < 1$. The parameter σ_{rr}^2 which is related to the RSI strength can be obtained at the pre-stage. Therefore, one straightforward choice of the power scaling factor is

$$\alpha^2 = \frac{P_r}{P_s + P_r\sigma_{rr}^2 + \sigma_n^2} \quad (4)$$

Since the m th tap of the overall IIR channel, which is θ^{m-1} , decreases in amplitude with increasing tap index m , we can assume that most of the energy (e.g. 99%) is contained in a finite length of the overall channel impulse response [10]. Define L as the effective length of the overall channel impulse response. Thus, we use a block-based transmission with a guard time of L symbol intervals to avoid inter-block interference [17]. At the receiver, it receives $N + L$ symbols and discards the last L symbols. Without loss of generality, the block length N is assumed to be far greater than L , so the rate loss due to the guard time is negligible. With the effective length L and block-based transmission, we can truncate the IIR channel. Let \mathbf{H}_θ be the matrix form of the channel in one block, which is given by an $N \times N$ Toeplitz matrix with first column $[1, \theta, \theta^2, \dots, \theta^{L-1}, 0, \dots, 0]^T$ and first row $[1, 0, \dots, 0]$.

We rewrite the output in terms of $\mathbf{x} := [x[0], \dots, x[N-1]]^T$ and $\mathbf{y} := [y_d[1], \dots, y_d[N]]^T$ as:

$$\mathbf{y} = h\mathbf{H}_\theta\mathbf{x} + d\mathbf{H}_\theta\mathbf{n}_r + \mathbf{n}_d, \quad (5)$$

where \mathbf{n}_r and \mathbf{n}_d are noise vectors composed of independent samples from the same distribution as $n_r[n]$ and $n_d[n]$ respectively. As can be seen from the matrix expression, $h\mathbf{H}_\theta$ is the overall channel and the sum of the last two terms in (5) is the colored noise. Thus, the overall channel becomes an ISI channel.

III. CHANNEL ESTIMATION

A. Maximum Likelihood Formulation

We now derive the ML estimator of h and θ for a given training sequence \mathbf{x} . We are only interested in h and θ since knowing them is enough for detection and equalization. So d is a nuisance parameter in our signal model. To derive the ML estimator, first we obtain the likelihood function. Given h , θ , and d , the mean and the covariance matrix of \mathbf{y} are

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{y}] = h\mathbf{H}_\theta\mathbf{x}, \quad (6)$$

$$\mathbf{C} = d^2\sigma_n^2\mathbf{H}_\theta\mathbf{H}_\theta^H + \sigma_n^2\mathbf{I}_N. \quad (7)$$

The likelihood function of \mathbf{y} is

$$p(\mathbf{y}; h, \theta, d) = \frac{1}{\pi^N |\mathbf{C}|} \exp \left(-(\mathbf{y} - \boldsymbol{\mu})^H \mathbf{C}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right), \quad (8)$$

where $|\mathbf{C}|$ denotes the determinant of matrix \mathbf{C} . The corresponding log-likelihood function is

$$\log p(\mathbf{y}; h, \theta, d) = -N \log \pi - \log |\mathbf{C}| - (\mathbf{y} - \boldsymbol{\mu})^H \mathbf{C}^{-1} (\mathbf{y} - \boldsymbol{\mu}). \quad (9)$$

Maximizing (9) is equivalent to minimizing the last two terms of it. Let f denote our objective function.

$$f(h, \theta, d) = \log |\mathbf{C}| + (\mathbf{y} - \boldsymbol{\mu})^H \mathbf{C}^{-1} (\mathbf{y} - \boldsymbol{\mu}). \quad (10)$$

The ML estimator is given by

$$\{\hat{h}, \hat{\theta}, \hat{d}\} = \arg \min_{h, \theta, d} \left\{ \log |\mathbf{C}| + (\mathbf{y} - \boldsymbol{\mu})^H \mathbf{C}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right\}. \quad (11)$$

Note that the three parameters are complex. We denote $h = h_x + jh_y$ where h_x and h_y are the real part and imaginary part of h respectively and j is the imaginary unit. Similarly we have $\theta = \theta_x + j\theta_y$ and $d = d_x + jd_y$. Before we solve the ML estimator, we will simplify the objective function to express it in terms of only one complex parameter θ . First we take derivative of f with

respect to h ,

$$\frac{\partial f}{\partial h} = -\mathbf{y}^H \mathbf{C}^{-1} \mathbf{H}_\theta \mathbf{x} + h^* \mathbf{x}^H \mathbf{H}_\theta^H \mathbf{C}^{-1} \mathbf{H}_\theta \mathbf{x}. \quad (12)$$

Setting the derivative to 0 we have

$$h = (\mathbf{x}^H \mathbf{H}_\theta^H \mathbf{C}^{-1} \mathbf{H}_\theta \mathbf{x})^{-1} \mathbf{x}^H \mathbf{H}_\theta^H \mathbf{C}^{-1} \mathbf{y}. \quad (13)$$

We can substitute (13) into (10) to eliminate h .

We now discuss how the nuisance parameter d affects the optimization of f . From (5), it is clear that parameter d is not identifiable since there is no training sequence incorporated with it. In what follows, we justify why the minimization in (11) with respect to d is unnecessary, and will not improve our estimates of the desired parameters, h and θ . In fact, we show in Appendix I that

$$\mathbb{E} \left[\frac{\partial f}{\partial d} \right] = 0. \quad (14)$$

Thus optimizing f with respect to d when other parameters are fixed does not affect the objective function on the average. Moreover, d has little impact on the optimization of h and θ since

$$\frac{\partial^2 f}{\partial h_x \partial d_x} \approx 0, \quad \frac{\partial^2 f}{\partial \theta_x \partial d_x} \approx 0, \quad (15)$$

as also shown in Appendix I. Since we optimize the complex parameters by optimizing their real parts and imaginary parts separately, we are able to only focus on $\frac{\partial^2 f}{\partial h_x \partial d_x}$ and $\frac{\partial^2 f}{\partial \theta_x \partial d_x}$. The derivatives of the imaginary parts are similar to (15). To summarize, h can be eliminated using (13), and d has little influence on f and therefore can be arbitrarily set. Therefore, (10) is a function only of one parameter, θ .

The objective function is not convex with respect to θ . To solve the problem numerically, we use the Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm [22], which is a popular quasi-Newton method. Note that the constraint $|\theta| < 1$ is imposed to ensure stability in (1). Euclidean projection is further applied to $\hat{\theta}$ to ensure that $|\hat{\theta}| < 1$ so that the estimates conform with the stability assumption. Because the algorithm can only deal with real valued parameters, the real and imaginary parts are optimized separately.

B. BFGS algorithm

We use the BFGS algorithm which is also used in [17] to solve the ML estimator in a different two-way relay context. The algorithm needs the gradients of f with respect to θ_x and θ_y . We derive the gradients in Appendix III. A linear MMSE estimator is used to initialize the BFGS algorithm, which helps the algorithm to converge faster and to reduce the possibility of trapping in a local minimum. We now elaborate on the initialization before we provide the details of the BFGS algorithm.

Pairs of received samples can be used for linear MMSE estimation even though the received samples \mathbf{y} are not linear in the desired parameters h and θ . We take two received symbols $y_d[2]$ and $y_d[3]$ to estimate h and θ . For estimating h , the second received symbol at the destination is used, which is

$$y_d[2] = hx[1] + dn_r[1] + n_d[2]. \quad (16)$$

The linear MMSE estimator \hat{h}_0 is given by

$$\hat{h}_0 = \frac{\alpha^2 \sigma_{sr}^2 \sigma_{rd}^2 x^*[1] y_d[2]}{\alpha^2 \sigma_{sr}^2 \sigma_{rd}^2 |x[1]|^2 + \alpha^2 \sigma_{rd}^2 \sigma_n^2 + \sigma_n^2}. \quad (17)$$

Let \tilde{h}_0 be the residual estimation error of h , i.e. $h = \hat{h}_0 + \tilde{h}_0$. Thus, \tilde{h}_0 is a random variable with zero mean and variance $\sigma_{\tilde{h}_0}^2$ which is given by

$$\sigma_{\tilde{h}_0}^2 = \frac{\alpha^2 \sigma_{sr}^2 \sigma_{rd}^2 (\alpha^2 \sigma_{rd}^2 \sigma_n^2 + \sigma_n^2)}{\alpha^2 \sigma_{sr}^2 \sigma_{rd}^2 |x[1]|^2 + \alpha^2 \sigma_{rd}^2 \sigma_n^2 + \sigma_n^2}. \quad (18)$$

After having \hat{h}_0 , we can estimate θ . First we use \hat{h}_0 to remove the known part $\hat{h}_0 x[2]$ in $y_d[3]$. The remaining signal of $y_d[3]$ is as follows:

$$y'_d[3] = \hat{h}_0 \theta x[1] + \tilde{h}_0 \theta x[1] + \tilde{h}_0 x[2] + h dn_r[1] + dn_r[2] + n_d[3]. \quad (19)$$

The linear MMSE estimator of θ is

$$\hat{\theta}_0 = \frac{\hat{h}_0 \alpha^2 \sigma_{rr}^2 x^*[1] y'_d[3]}{\hat{h}_0^2 \alpha^2 \sigma_{rr}^2 |x[1]|^2 + \sigma_{\tilde{h}_0}^2 |x[2]|^2 + \sigma_{\tilde{h}_0}^2 \alpha^2 \sigma_{rr}^2 |x[1]|^2 + \alpha^4 \sigma_{sr}^2 \sigma_{rd}^4 \sigma_n^2 + \alpha^2 \sigma_{rd}^2 \sigma_n^2 + \sigma_n^2}. \quad (20)$$

Though we only make use of one training symbol in the above linear MMSE method, it is possible to extend the method to use multiple symbols, in which case a special training sequence with $L - 1$ zeros followed by one symbol is transmitted, where L is the effective length of the channel impulse response.

We now provide the BFGS algorithm which uses the initialization explained above, and the gradients in Appendix III.

Initialize: $\mathbf{z}_0 \triangleq [\hat{\theta}_x \ \hat{\theta}_y]^T$, $\mathbf{A}_0^{-1} = \mathbf{I}_{2 \times 2}$.

Repeat until convergence for k : (BFGS)

1. Obtain a search direction $\mathbf{p}_k = -\mathbf{A}_k^{-1} \nabla f(\mathbf{z}_k)$.
2. Find stepsize λ_k by backtracking linesearch, then update $\mathbf{z}_{k+1} = \mathbf{z}_k + \lambda_k \mathbf{p}_k$.
3. Set $\mathbf{s}_k = \lambda_k \mathbf{p}_k$, $\mathbf{v}_k = \nabla f(\mathbf{z}_{k+1}) - \nabla f(\mathbf{z}_k)$
4. Update the inverse Hessian approximation by

$$\mathbf{A}_{k+1}^{-1} = \mathbf{A}_k^{-1} + \frac{(\mathbf{s}_k^T \mathbf{v}_k + \mathbf{v}_k^T \mathbf{A}_k^{-1} \mathbf{v}_k) \mathbf{s}_k \mathbf{s}_k^T}{(\mathbf{s}_k^T \mathbf{v}_k)^2} - \frac{\mathbf{A}_k^{-1} \mathbf{v}_k \mathbf{s}_k^T + \mathbf{s}_k \mathbf{v}_k^T \mathbf{A}_k^{-1}}{\mathbf{s}_k^T \mathbf{v}_k}$$

Obtain the converged result \mathbf{z}_k .

If $|\hat{\theta}| > 1$, $\hat{\theta} = \hat{\theta}/|\hat{\theta}|$ (Euclidean projection).

After the initial values are input, θ is optimized by the BFGS algorithm. The iteration is controlled by the index k . The results of one iteration will be used as initial values for the next iteration. In particular, there is a constraint $|\theta| < 1$ on θ . We use Euclidean projection to keep $\hat{\theta}$ in its valid region. If the result of $\hat{\theta}$ is a point outside of the valid region, Euclidean projection maps the outside point to its nearest valid point.

It is guaranteed that the BFGS algorithm converges to a local minimum point because it is a descent algorithm [22, Sec. 8.3.5]. However, for non-convex objective functions, the convergence point may be a local minimum. To avoid this, we use MMSE estimates of the parameters to initialize the algorithm as mentioned above.

The complexity of the algorithm is dominated by the matrix inversion of the covariance matrix \mathbf{C} in the calculation of the gradients and the objective function (10). For large training length N ,

\mathbf{C} asymptotically becomes to a positive definite Toeplitz matrix. The complexity of inverting it is $O(N \log^2 N)$. The approximate inverse-Hessian matrix update only depends on the number of parameters to be estimate but not on N . Therefore, the total complexity of the BFGS algorithm in one iteration is $O(N \log^2 N)$ for large N . Moreover, the algorithm with linear MMSE initialization converges faster than random initialization based on our observation in the simulation. Thus, our initialization method also helps to reduce the complexity of the algorithm.

IV. OPTIMAL TRAINING SEQUENCES

A. Cramer-Rao Bounds

The CRB is derived not only to show the accuracy of the estimates but also to act as a metric when designing the training sequences. Differentiating the log-likelihood function $\log p(\mathbf{y}; h, \theta, d)$ twice, we can obtain the Fisher information matrix (FIM). We are only interested in h and θ and set d as a nuisance parameter. In our method, the CRBs we derived are the CRBs conditioned on d . Let $\boldsymbol{\xi} = [h \ \theta]^T$ be the vector of parameters. The FIM is given by

$$\mathbf{\Gamma}(\boldsymbol{\xi}) = \mathbb{E} \left[\frac{\partial \log p}{\partial \boldsymbol{\xi}^*} \frac{\partial \log p}{\partial \boldsymbol{\xi}^T} \right]. \quad (21)$$

The (m, n) element of $\mathbf{\Gamma}$ is given by

$$\Gamma_{mn} = \frac{\partial \boldsymbol{\mu}^H}{\partial \xi_m^*} \mathbf{C}^{-1} \frac{\partial \boldsymbol{\mu}}{\partial \xi_n} + \text{tr} \left(\mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \xi_m^*} \mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \xi_n} \right), \quad (22)$$

where ξ_m is the m th element of $\boldsymbol{\xi}$. Thus we have

$$\Gamma_{11} = \frac{\partial \boldsymbol{\mu}^H}{\partial h^*} \mathbf{C}^{-1} \frac{\partial \boldsymbol{\mu}}{\partial h} + \text{tr} \left(\mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial h^*} \mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial h} \right) = \mathbf{x}^H \mathbf{H}_\theta^H \mathbf{C}^{-1} \mathbf{H}_\theta \mathbf{x}. \quad (23)$$

Similarly,

$$\begin{aligned} \Gamma_{22} &= \frac{\partial \boldsymbol{\mu}^H}{\partial \theta^*} \mathbf{C}^{-1} \frac{\partial \boldsymbol{\mu}}{\partial \theta} + \text{tr} \left(\mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta^*} \mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta} \right) \\ &= |h|^2 \mathbf{x}^H \mathbf{B}_\theta^H \mathbf{C}^{-1} \mathbf{B}_\theta \mathbf{x} + |d|^4 \sigma_n^4 \text{tr} \left(\mathbf{C}^{-1} \mathbf{H}_\theta \mathbf{B}_\theta^H \mathbf{C}^{-1} \mathbf{B}_\theta \mathbf{H}_\theta^H \right). \end{aligned} \quad (24)$$

Since \mathbf{C} is not a function of h ,

$$\Gamma_{12} = \frac{\partial \boldsymbol{\mu}^H}{\partial h^*} \mathbf{C}^{-1} \frac{\partial \boldsymbol{\mu}}{\partial \theta} = h \mathbf{x}^H \mathbf{H}_\theta^H \mathbf{C}^{-1} \mathbf{B}_\theta \mathbf{x}, \quad (25)$$

$$\Gamma_{21} = \frac{\partial \boldsymbol{\mu}^H}{\partial \theta^*} \mathbf{C}^{-1} \frac{\partial \boldsymbol{\mu}}{\partial h} = h^* \mathbf{x}^H \mathbf{B}_\theta^H \mathbf{C}^{-1} \mathbf{H}_\theta \mathbf{x}. \quad (26)$$

The CRB is given by the trace of the inverse of $\boldsymbol{\Gamma}$, which is $CRB_\xi = \text{tr}(\boldsymbol{\Gamma}^{-1})$. In particular, the CRBs for each parameter are the diagonal elements of the inverse FIM. Since $\boldsymbol{\Gamma}$ is a 2 by 2 complex matrix, we can find its inverse by calculating its determinant and adjoint. The determinant is $|\boldsymbol{\Gamma}| = \Gamma_{11}\Gamma_{22} - \Gamma_{12}\Gamma_{21}$. Therefore, the CRBs are given by $CRB_h = \Gamma_{22}/|\boldsymbol{\Gamma}|$, $CRB_\theta = \Gamma_{11}/|\boldsymbol{\Gamma}|$.

B. Training Sequence Design via the CRB

In this subsection, we analyze the CRBs by using Szegő's theorem [23], [24] and minimize it in the regime where the training length N is large. We show that the optimal training sequence that minimizes the CRB is sinusoidal and we characterize the frequency of this sinusoidal. To do the asymptotic behavior of the Toeplitz matrix, \mathbf{H}_θ will be used. Define a function

$$t(\lambda) = \sum_{k=0}^{\infty} t_k e^{j\lambda k}, \quad (27)$$

where t_k are the elements of the first column of the $N \times N$ Toeplitz matrix \mathbf{T}_N . Thus we can use the function $t(\lambda)$ to represent the matrix. We denote the matrix as $\mathbf{T}_N(t(\lambda))$. In our system, $t_k = \theta^k$ for $k = 0, \dots, L-1$ and otherwise $t_k = 0$. We show that \mathbf{H}_θ and $\mathbf{T}_N(t(\lambda))$ are asymptotically equivalent. First, since both \mathbf{H}_θ and $\mathbf{T}_N(t(\lambda))$ are banded Toeplitz matrices, their strong norms (operator norms) are bounded. Secondly, from the definition of $\mathbf{T}_N(t(\lambda))$, we have $\lim_{N \rightarrow \infty} \|\mathbf{H}_\theta - \mathbf{T}_N(t(\lambda))\| = 0$, where $\|\mathbf{A}\|$ denotes the weak norm (Hilbert-Schmidt norm) of matrix \mathbf{A} . With the two conditions above, we can say that \mathbf{H}_θ and $\mathbf{T}_N(t(\lambda))$ are asymptotically equivalent [24]. Therefore, we will write $\mathbf{H}_\theta = \mathbf{T}_N(t(\lambda))$ which will be understood to hold for asymptotically large N . We have the following expression for $t(\lambda)$,

$$t(\lambda) = \sum_{k=0}^{L-1} \theta^k e^{j\lambda k} = \frac{1 - \theta^L e^{j\lambda L}}{1 - \theta e^{j\lambda}} = \frac{1}{1 - \theta e^{j\lambda}}, \quad (28)$$

where $\theta^L \approx 0$ by our assumption of channel energy in Section II. The covariance matrix \mathbf{C} is

$$\mathbf{C} = |d|^2 \sigma_n^2 \mathbf{T}_N(t(\lambda)) \mathbf{T}_N(t^*(\lambda)) + \sigma_n^2 \mathbf{I}_N. \quad (29)$$

Without loss of generality, we set $\sigma_n^2 = 1$. According to Szegő's theorem, the product of two Toeplitz matrices is a Toeplitz matrix asymptotically, as well as the inverse of a Toeplitz matrix. Thus we have

$$\mathbf{C} = |d|^2 \mathbf{T}_N(|t(\lambda)|^2) + \mathbf{I}_N = \mathbf{T}_N(|d|^2 |t(\lambda)|^2 + 1), \quad (30)$$

$$\mathbf{C}^{-1} = \mathbf{T}_N\left(\frac{1}{|d|^2 |t(\lambda)|^2 + 1}\right). \quad (31)$$

The Fisher information of the source-relay-destination channel h is

$$\Gamma_{11} = \mathbf{x}^H \mathbf{H}_\theta^H \mathbf{C}^{-1} \mathbf{H}_\theta \mathbf{x} = \mathbf{x}^H \mathbf{T}_N(t^*(\lambda)) \mathbf{T}_N\left(\frac{1}{|d|^2 |t(\lambda)|^2 + 1}\right) \mathbf{T}_N(t(\lambda)) \mathbf{x} \quad (32)$$

$$= \frac{|t(\lambda)|^2}{|d|^2 |t(\lambda)|^2 + 1} \|\mathbf{x}\|^2, \quad (33)$$

where $\frac{|t(\lambda)|^2}{|d|^2 |t(\lambda)|^2 + 1}$ is the eigenvalue of $\mathbf{T}_N\left(\frac{|t(\lambda)|^2}{|d|^2 |t(\lambda)|^2 + 1}\right)$ and depends on λ .

Similarly, we can denote $\mathbf{B}_\theta = \mathbf{T}_N(g(\lambda))$ where $g(\lambda) = \frac{e^{j\lambda}}{(1 - \theta e^{j\lambda})^2}$ is the derivative of $t(\lambda)$ with respect to θ . Then the Fisher information of the RSI channel can be represented by Toeplitz matrices as

$$\Gamma_{22} = |h|^2 \mathbf{x}^H \mathbf{B}_\theta^H \mathbf{C}^{-1} \mathbf{B}_\theta \mathbf{x} + |d|^4 \text{tr}(\mathbf{C}^{-1} \mathbf{H}_\theta \mathbf{B}_\theta^H \mathbf{C}^{-1} \mathbf{B}_\theta \mathbf{H}_\theta^H) \quad (34)$$

$$= |h|^2 \mathbf{x}^H \mathbf{T}_N\left(\frac{|g(\lambda)|^2}{|d|^2 |t(\lambda)|^2 + 1}\right) \mathbf{x} + |d|^4 \text{tr}\left(\mathbf{T}_N\left(\frac{|g(\lambda)|^2 |t(\lambda)|^2}{(|d|^2 |t(\lambda)|^2 + 1)^2}\right)\right) \quad (35)$$

$$= |h|^2 \frac{|g(\lambda)|^2}{|d|^2 |t(\lambda)|^2 + 1} \|\mathbf{x}\|^2 + |d|^4 \frac{\|\mathbf{x}\|^2}{2\pi P_s} \int_0^{2\pi} \frac{|g(\lambda)|^2 |t(\lambda)|^2}{(|d|^2 |t(\lambda)|^2 + 1)^2} d\lambda. \quad (36)$$

We can simplify the first term in (36) similarly to Γ_{11} . The second term comes from the fact that the trace of Toeplitz matrices is equal to the integral of the function of λ that characterizes it. Simplifying this integral we have

$$\int_0^{2\pi} \frac{|g(\lambda)|^2 |t(\lambda)|^2}{(|d|^2 |t(\lambda)|^2 + 1)^2} d\lambda = \int_0^{2\pi} \frac{1}{|1 - \theta e^{j\lambda}|^2 (|d|^2 + |1 - \theta e^{j\lambda}|^2)^2} d\lambda. \quad (37)$$

In our FD relay system, we assume $P_r \gg P_s$ since P_s is the received power at the relay which incorporates the pathloss. Thus $|d|^2 \gg 1$. Note that $|\theta| < 1$, we have $|d|^2 \gg |1 - \theta e^{j\lambda}|^2$. We can approximate the integral as

$$\begin{aligned} & \frac{1}{2\pi} \int_0^{2\pi} \frac{1}{|1 - \theta e^{j\lambda}|^2 (|d|^2 + |1 - \theta e^{j\lambda}|^2)^2} d\lambda \\ & \approx \frac{1}{2\pi} \int_0^{2\pi} \frac{1}{|1 - \theta e^{j\lambda}|^2 |d|^4} d\lambda = \frac{1}{|d|^4} \frac{1}{(|\theta| + 1)|\theta| - 1}. \end{aligned} \quad (38)$$

Therefore, the Fisher information of the RSI channel θ becomes

$$\Gamma_{22} = \frac{|g(\lambda)|^2}{|d|^2 |t(\lambda)|^2 + 1} \|\mathbf{x}\|^2 + \frac{1}{P_s(|\theta| + 1)|\theta| - 1} \|\mathbf{x}\|^2. \quad (39)$$

Similarly, we can represent Γ_{12} and Γ_{21} as

$$\Gamma_{12} = h \mathbf{x}^H \mathbf{B}_\theta^H \mathbf{C}^{-1} \mathbf{H}_\theta \mathbf{x} = \mathbf{x}^H \mathbf{T}_N \left(\frac{ht(\lambda)g^*(\lambda)}{|d|^2 |t(\lambda)|^2 + 1} \right) \mathbf{x}, \quad (40)$$

$$\Gamma_{21} = h^* \mathbf{x}^H \mathbf{H}_\theta^H \mathbf{C}^{-1} \mathbf{B}_\theta \mathbf{x} = \mathbf{x}^H \mathbf{T}_N \left(\frac{h^* t^*(\lambda)g(\lambda)}{|d|^2 |t(\lambda)|^2 + 1} \right) \mathbf{x}. \quad (41)$$

To calculate the CRB, we need the product of Γ_{12} and Γ_{21} which is

$$\Gamma_{12}\Gamma_{21} = |h|^2 \frac{p^2(\lambda)}{(|d|^2 |t(\lambda)|^2 + 1)^2} \|\mathbf{x}\|^4. \quad (42)$$

where $p(\lambda) = \frac{1}{2}[t^*(\lambda)g(\lambda) + t(\lambda)g^*(\lambda)]$. Function $p(\lambda)$ is the real part of $t^*(\lambda)g(\lambda)$ and it shows that only the symmetric part of the Toeplitz matrix affects the product.

Thus, the CRB of θ is

$$CRB_\theta = \frac{\Gamma_{11}}{\Gamma_{11}\Gamma_{22} - \Gamma_{12}\Gamma_{21}} \quad (43)$$

$$= \frac{1}{\|\mathbf{x}\|^2} \frac{|t(\lambda)|^2 (\alpha^2 |t(\lambda)|^2 + 1)}{|h|^2 |t(\lambda)|^2 |g(\lambda)|^2 + A |t(\lambda)|^2 (|d|^2 |t(\lambda)|^2 + 1) - |h|^2 |p(\lambda)|^2} \quad (44)$$

$$\triangleq \frac{1}{\|\mathbf{x}\|^2} F(\lambda), \quad (45)$$

where $A = (P_s(|\theta| + 1)|\theta| - 1)^{-1}$. To find the frequency of the sinusoidal training sequence, we

minimize the CRB in (45):

$$\min_{\lambda \in [0, 2\pi]} F(\lambda)$$

Simplifying $F(\lambda)$ we have

$$F(\lambda) = \frac{\alpha^2 + |1 - \theta e^{j\lambda}|^2}{\frac{|h|^2}{|1 - \theta e^{j\lambda}|^2} + A(\alpha^2 + |1 - \theta e^{j\lambda}|^2) - \frac{1}{4}|h|^2 \frac{(\text{Re}[e^{j\lambda} - \theta^*])^2}{|1 - \theta e^{j\lambda}|^4}}. \quad (46)$$

Let $x = |1 - \theta e^{j\lambda}|^2$ and $x \in [(1 - |\theta|^2), (1 + |\theta|^2)]$. Substitute it into $F(\lambda)$, we have

$$G(x) = \frac{\alpha^2 + x}{\frac{|h|^2}{x} + A(\alpha^2 + x) - \frac{|h|^2}{4x^2} \left(\frac{(1 + |\theta|^2 - x)\theta_R}{2|\theta|^2} + \sqrt{1 - \frac{(1 + |\theta|^2 - x)^2}{4|\theta|^2}} \frac{\theta_I}{|\theta|} - \theta_R \right)^2}, \quad (47)$$

where θ_R and θ_I are the real and imaginary parts of θ respectively.

The optimal solution that minimizes $F(\lambda)$ can be found by numerically solving $G'(x) = 0$. $G'(x)$ is a polynomial in x with the highest order 8. The coefficients of the polynomial are given in Appendix II. Note that $x \in [(1 - |\theta|^2), (1 + |\theta|^2)]$, so that the two endpoints of the interval are also candidates for the optimal solution in case that the only solution to $G'(x) = 0$ is a saddle point or there is no solution in the interval. After we get all the candidates (of which there are a maximum of 8), we are able to substitute each of them into $G(x)$ to find the one that minimizes the function. Assume the solution that minimizes CRB_θ is λ^* , the corresponding optimal training sequence is given by $\frac{1}{\sqrt{N}}[1, \dots, e^{j2\pi k\lambda^*}, \dots, e^{j2\pi(N-1)\lambda^*}]^T$ for $k = 0, \dots, N - 1$.

The CRB for h also can be derived the same way as the CRB of θ , and also can be minimized by finding the roots of a polynomial. We do not include it here due to lack of space. We calculate the CRB of h in the simulations. When we minimize the CRB for θ , it is not guaranteed that the CRB of h is minimized as well. However, we can minimize the sum of CRBs of θ and h if both parameters are considered, also through polynomial rooting.

The optimal training sequence depends on both of the channel h and θ through λ . In practice, we do not have the information of h and θ until the first training sequence is sent. We can apply an adaptive training method where the optimal training sequence is designed by using estimates obtained from its previous training sequence. In what follows we show through an approximation

that the minimizer of (47) only weakly depends on h .

C. Low Complexity Approximation

The optimal solution can be found by minimizing the CRB numerically via finding the polynomial roots. However, the complexity can be reduced by a certain approximation which we now describe. This provides an approximately optimal and practical solution for the problem. Assume $|\theta|$ is small, so that the value of x is very close to 1. Then we can have the following approximation

$$\left(\frac{(1 + |\theta|^2 - x)\theta_R}{2|\theta|^2} + \sqrt{1 - \frac{(1 + |\theta|^2 - x)^2}{4|\theta|^2}} \frac{\theta_I}{|\theta|} - \theta_R \right)^2 \approx 1. \quad (48)$$

Thus,

$$G(x) \approx \frac{\alpha^2 + x}{\frac{|h|^2}{x} + A(\alpha^2 + x) - \frac{|h^2|}{4x^2}}. \quad (49)$$

Solving $G'(x) = 0$ is equivalent to solving the following equation:

$$8|h^2|x^3 + (4\alpha^2 - 3)|h|^2x^2 - 2\alpha^2|h|^2x = 0. \quad (50)$$

Equation (50) shows that h does not affect the solution of $G'(x) = 0$. One can verify that none of the three real roots of (50) is in the valid interval of x which is $[(1 - |\theta|)^2, (1 + |\theta|)^2]$. Note that $G(x)$ is an increasing function since $|h|^2 > 0$. Therefore the approximately optimal solution is the left endpoint of the interval i.e. $x = (1 - |\theta|)^2$. Thus $\lambda = -\angle\theta$ where \angle represents the phase of a complex number. Moreover, the channel $|h|$ does not affect the solution of x , so that the training sequence for estimating θ only depends on θ and not $|h|$, making it easier to implement than the optimal training sequence. The normalized corresponding training sequence is given by $\frac{1}{\sqrt{N}}[1, \dots, e^{j2\pi k\lambda_1^*}, \dots, e^{j2\pi(N-1)\lambda_1^*}]^T$ for $k = 0, \dots, N - 1$ where λ_1^* minimizes (49).

V. FREQUENCY SELECTIVE CHANNELS

In this section, we extend our channel estimation method to the case where the channels between nodes are frequency selective fading. We show that our training method and Fisher information calculation can be extended to this case.

We assume the source-to-relay and relay-to-destination channels are frequency selective fading with channel taps $\mathbf{h}_{\text{sr}} = [h_{\text{sr}}[1], h_{\text{sr}}[2], \dots, h_{\text{sr}}[L_1]]$ and $\mathbf{h}_{\text{rd}} = [h_{\text{rd}}[1], h_{\text{rd}}[2], \dots, h_{\text{rd}}[L_2]]$ respectively. Therefore, with block based transmission, the channel matrix for the source-to-relay channel \mathbf{H}_{sr} is an $N \times N$ Toeplitz matrix with first column $[\mathbf{h}_{\text{sr}}^T, 0, \dots, 0]^T$ and first row $[1, 0, 0, \dots, 0]$. For the relay-to-destination channel, the channel matrix \mathbf{H}_{rd} is also an $N \times N$ Toeplitz matrix with first column $[\mathbf{h}_{\text{rd}}^T, 0, \dots, 0]^T$ and first row $[1, 0, 0, \dots, 0]$. The received signal for the training phase is similar to (5) and becomes

$$\mathbf{y}_f = \alpha_f \mathbf{H}_{\text{rd}} \mathbf{H}_\theta \mathbf{H}_{\text{sr}} \mathbf{x} + \alpha_f \mathbf{H}_{\text{rd}} \mathbf{H}_\theta \mathbf{n}_r + \mathbf{n}_d, \quad (51)$$

where α_f is the new power scaling factor for the frequency selective channel given by

$$\alpha_f^2 = \frac{P_r}{P_s(1 + L_1 \sigma_n^2) + P_r \sigma_{\text{rr}}^2 + \sigma_n^2}. \quad (52)$$

The overall channel is $\alpha_f \mathbf{H}_{\text{rd}} \mathbf{H}_\theta \mathbf{H}_{\text{sr}}$ in the frequency selective case instead of $h \mathbf{H}_\theta$ for flat fading.

To extend our training method, we can use Szegő's theorem for Toeplitz matrices which is explained in Section IV-B to approximate the overall channel matrix as a Toeplitz matrix. According to the theorem, when the training length is large, the product of two Toeplitz matrices is still a Toeplitz matrix and the elements of the product can be determined by the elements of the two matrices. Similar to the way we define $\mathbf{H}_\theta = \mathbf{T}_N(t(\lambda))$ for large N , we can define $\mathbf{H}_{\text{sr}} = \mathbf{T}_N(q(\lambda))$ and $\mathbf{H}_{\text{rd}} = \mathbf{T}_N(p(\lambda))$, where $q(\lambda) = \sum_{k=0}^{L_1-1} h_{\text{sr}}[k+1] e^{j\lambda k}$ and $p(\lambda) = \sum_{k=0}^{L_2-1} h_{\text{rd}}[k+1] e^{j\lambda k}$. Thus, the overall channel matrix is

$$\mathbf{H}_f = \mathbf{H}_{\text{rd}} \mathbf{H}_\theta \mathbf{H}_{\text{sr}} = \alpha_f \mathbf{T}_N(p(\lambda)) \mathbf{T}_N(t(\lambda)) \mathbf{T}_N(q(\lambda)) \quad (53)$$

$$= \alpha_f \mathbf{T}_N(p(\lambda) t(\lambda) q(\lambda)), \quad (54)$$

for large N .

\mathbf{H}_f is also a Toeplitz matrix. Assume the elements in its first column are $h_f[k]$ for $k = 1, 2, \dots, N$,

we have

$$h_f[k] = \alpha_f \frac{1}{2\pi} \int_0^{2\pi} p(\lambda) t(\lambda) q(\lambda) e^{-jk\lambda} d\lambda. \quad (55)$$

The parameters to be estimated are $h_f[k]$ for $k = 1, 2, \dots, L_f$ where $L_f = L_1 + L_2 + L - 2$. Assume $\boldsymbol{\xi}_f = [h_f[1], \dots, h_f[L_f]]^T$, our ML estimator can be extended to estimate $\boldsymbol{\xi}_f$ by the following. First $h_f[1]$ has the same position as h in (11). Then using $h_f[i + 1]$ replace θ^i in (11). Thus, our ML method can be applied to estimate the overall channel even in the frequency selective setup.

The Fisher information for the frequency selective can be obtained similarly to the flat fading case by using

$$\Gamma_{mn}^{(f)} = \frac{\partial \boldsymbol{\mu}_f^H}{\partial \xi_{fm}^*} \mathbf{C}_f^{-1} \frac{\partial \boldsymbol{\mu}_f}{\partial \xi_{fn}} + \text{tr} \left(\mathbf{C}_f^{-1} \frac{\partial \mathbf{C}_f}{\partial \xi_{fm}^*} \mathbf{C}_f^{-1} \frac{\partial \mathbf{C}_f}{\partial \xi_{fn}} \right), \quad (56)$$

where

$$\boldsymbol{\mu}_f = \mathbf{H}_f \mathbf{x}, \quad \mathbf{C}_f = \alpha_f^2 \sigma_n^2 \mathbf{H}_{rd} \mathbf{H}_\theta \mathbf{H}_\theta^H \mathbf{H}_{rd}^H + \sigma_n^2 \mathbf{I}_N. \quad (57)$$

The CRBs are given by the diagonal elements of the inverse of the Fisher information matrix $\Gamma^{(f)}$. If desired, a single scalar quantity representing the overall CRB can be computed by find the trace of this matrix:

$$CRB_{\boldsymbol{\xi}_f} = \text{tr} \left((\Gamma^{(f)})^{-1} \right). \quad (58)$$

VI. MULTIPLE RELAYS

We now extend our training method to the multi-relay case which has fixed distance and equally-spaced relays between the source and the destination. Assume there are M relays between the source and the destination which satisfy $M(L - 1) < N$ (There is a guard time of $L - 1$ symbols for each relay). Each relay works in FD mode with AF relay protocol. The relays have their own RSI, and they do not perform estimation or equalization to keep relay complexity low. The estimation and equalization are performed only at the destination. The channel between the $(i - 1)$ th relay and the i th relay is flat fading with coefficient h_i for $i = 2, 3, \dots, M$. The channels from the source to the first relay and from the last relay to the destination are h_1 and h_{M+1} . Channel coefficients h_i for

$i = 1, \dots, M + 1$ are Gaussian random variables with zero-mean and variance σ_h^2 . Each relay has its own RSI channel h_{rri} and power scaling factor α_i which is given by

$$\alpha_i^2 = \frac{P_r}{P_{si}\sigma_h^2 + P_r\sigma_{\text{rr}}^2 + \sigma_n^2}, \quad (59)$$

where P_{si} is the received power of the desired signal at the i th relay. We also assume all the relays have the same average transmit power and average RSI power for simplicity.

The distance between the source and the destination is fixed in our model and M relays are placed in the line between the source and the destination in an equally spaced manner. Assume the distance between the source and the destination is normalized and the corresponding path loss is K dB. Then by using a simplified path loss model [25], the path loss between two relays is

$$K_m \text{dB} = K \text{dB} + 10\gamma \log_{10}(M + 1), \quad (60)$$

where γ is the path loss exponent. We incorporate the path loss into h_i which leads to $\sigma_h^2 = 1/K_m$ and $P_{si} = P_r K(M + 1)^\gamma$.

The transmit signal at the m th relay is

$$\mathbf{y}_m = \prod_{i=1}^m (\alpha_i h_i \mathbf{H}_{\theta_i}) \mathbf{x} + \sum_{i=1}^m \left[\left(\prod_{n=i+1}^m \alpha_n h_n \prod_{n=i}^m \mathbf{H}_{\theta_n} \right) \mathbf{n}_{\text{ri}} \right], \quad (61)$$

where \mathbf{H}_{θ_i} is the RSI channel at the i th relay defined similarly as \mathbf{H}_θ with $\theta_i = \alpha_i h_{\text{rri}}$, and \mathbf{n}_{ri} is the additive Gaussian white noise at the i th relay. The received signal at the destination from the m th relay is given by

$$\begin{aligned} \mathbf{y}_d &= h_{M+1} \mathbf{y}_M + \mathbf{n}_d \\ &= h_{M+1} \prod_{i=1}^M (\alpha_i h_i \mathbf{H}_{\theta_i}) \mathbf{x} + \sum_{i=1}^M \left[\left(\prod_{n=i+1}^M \alpha_n h_n \prod_{n=i}^M \mathbf{H}_{\theta_n} \right) \mathbf{n}_{\text{ri}} \right] + \mathbf{n}_d. \end{aligned} \quad (62)$$

Define $\mathbf{H}^{(n)} = \prod_{i=n}^M \mathbf{H}_{\theta_i}$ and its corresponding function $t^{(n)}(\lambda)$. By using the property of product

of Toeplitz matrix for large N , we have

$$\mathbf{H}^{(n)} = \prod_{i=n}^M \mathbf{T}(t_{\theta_i}(\lambda)) = \mathbf{T}\left(\prod_{i=n}^M t_{\theta_i}(\lambda)\right), \quad (63)$$

where $t_{\theta_i}(\lambda)$ is defined the same as (28) with θ_i . $\mathbf{H}^{(n)}$ is also a Toeplitz matrix and the elements in its first row defined through an inverse Fourier transform

$$h_k^{(n)} = \frac{1}{2\pi} \int_0^{2\pi} t^{(n)}(\lambda) e^{-jk\lambda} d\lambda, \quad (64)$$

where $t^{(n)}(\lambda) = \prod_{i=n}^M t_{\theta_i}(\lambda)$. Thus we can approximate \mathbf{y}_d for large N as

$$\mathbf{y}_d = z_M \mathbf{H}^{(1)} \mathbf{x} + \sum_{i=1}^M \left[\left(\prod_{n=i}^M \alpha_n h_n \right) / (\alpha_i h_i) \mathbf{H}^{(i)} \mathbf{n}_{ri} \right] + \mathbf{n}_d, \quad (65)$$

where $z_M = h_{M+1} \prod_{i=1}^M \alpha_i h_i$. The first row of $\mathbf{H}^{(1)}$ is $[1, h_2^{(1)}, h_3^{(1)}, \dots, h_{M(L-1)}^{(1)}, 0, \dots, 0]^T$ which is an $N \times 1$ vector (Assume $M(L-1) < N$). The channel parameters to be estimated is $\boldsymbol{\xi}_M = [z_M, h_2^{(1)}, h_3^{(1)}, \dots, h_{M(L-1)}^{(1)}]^T$. The signal model of (65) is the same to that of (51) except additional noise terms. We can also extend our ML estimator to estimate the channel parameters using the same way in the frequency selective case where z_M and $h_k^{(1)}$ are analogous to h and θ^k in (51) respectively. The advantage of estimating the multi-relay channel at the destination rather than at each relay is to keep the relays low complexity with just analog signal processing capability. However, the performance is better if estimation and equalization are performed at each relay, at the cost of complexity.

The CRB for multiple relays can be derived and analyzed similarly to the single relay case. We derive the CRBs for the first two strongest channel taps z_M and $h_2^{(1)}$ which dominate the data detection. The CRBs for other parameters can also be found similarly to (56). The CRB of $h_2^{(1)}$ is given by

$$CRB_{h_2^{(1)}} = \frac{|t^{(M)}(\lambda)|^2 (\sum_{i=1}^M |c_i|^2 |t_{\theta_i}(\lambda)|^2 + 1)}{J \left(\sum_{i=1}^M |c_i|^2 |t_{\theta_i}(\lambda)|^2 + 1 \right) + \|\mathbf{x}\|^2 \|z_M\|^2 (|t^{(M)}(\lambda)|^2 |g^{(M)}(\lambda)|^2 - |p^{(M)}(\lambda)|^2)}, \quad (66)$$

where $c_i = \prod_{n=i}^M \alpha_n h_n / (\alpha_i h_i)$ and

$$J = |c_1|^4 \frac{1}{2\pi} \int_0^{2\pi} \frac{|t^{(M)}(\lambda)|^2 |g^{(M)}(\lambda)|^2}{(\sum_{i=1}^M |c_i|^2 |t_{\theta_i}(\lambda)|^2 + 1)^2}. \quad (67)$$

Define functions $g^{(M)}(\lambda) = \frac{\partial t^{(M)}(\lambda)}{\partial \lambda}$ and $p^{(M)}(\lambda) = \frac{1}{2} [(t^{(M)}(\lambda))^* g^{(M)}(\lambda) + t^{(M)}(\lambda) (g^{(M)}(\lambda))^*]$. For z_M we have

$$CRB_{z_M} = \frac{[|z_M|^2 |g^{(M)}(\lambda)|^2 + J(\sum_{i=1}^M |c_i|^2 |t_{\theta_i}(\lambda)|^2 + 1)](\sum_{i=1}^M |c_i|^2 |t_{\theta_i}(\lambda)|^2 + 1)}{J \left(\sum_{i=1}^M |c_i|^2 |t_{\theta_i}(\lambda)|^2 + 1 \right) + \|\mathbf{x}\|^2 |z_M|^2 (|t^{(M)}(\lambda)|^2 |g^{(M)}(\lambda)|^2 - |p^{(M)}(\lambda)|^2)}. \quad (68)$$

We further approximate the CRBs and have simple expressions to find how the number of relays affects the CRBs. From (27) we have $|t_{\theta_i}(\lambda)| \approx 1$ for small θ_i . We also assume the training power is large so that $J \ll \|\mathbf{x}\|^2$. Therefore the CRB becomes

$$CRB_{h_2^{(1)}} \approx \frac{\sum_{i=1}^M |c_i|^2 + 1}{|z_M|^2 \|\mathbf{x}\|^2}, \quad (69)$$

$$CRB_{z_M} \approx \frac{\sum_{i=1}^M |c_i|^2 + 1}{\|\mathbf{x}\|^2}. \quad (70)$$

By plugging in $|c_i|^2 = (\alpha_i^2)^{M-i} \prod_{n=i+1}^M |h_i|^2$, and $\alpha_i^2 = \frac{P_r}{P_r K(M+1)^\gamma + P_r \sigma_{rr}^2 + 1}$, we have the CRB expressions as a function of M .

$$CRB_{h_{M2}} = \frac{\sum_{i=1}^M (K(M+1)^\gamma + k_1)^{i-M} (\prod_{n=i+1}^M |h_i|^2)^{i-M} + 1}{\|\mathbf{x}\|^2 (L(M+1)^\gamma + k_1)^{-M} \prod_{n=i+1}^M |h_i|^2}, \quad (71)$$

$$CRB_{z_M} = \frac{\sum_{i=1}^M (K(M+1)^\gamma + k_1)^{i-M} (\prod_{n=i+1}^M |h_i|^2)^{i-M} + 1}{\|\mathbf{x}\|^2}. \quad (72)$$

where $k_1 = \sigma_{rr}^2 + 1/P_r$. Equation (71) and (72) are simple functions of M . Intuitively, the estimates of z_M will become more inaccurate as the noise goes strong for increasing M . However, as the number of relays increases, the RSI for each relay accumulates at the destination, which makes the RSI channel h_{M2} stronger and easier to estimate. Thus there is an optimal M which minimizes the sum MSE of z_M and h_{M2} . Since M is an integer and is not quite large, the optimal number of relays with respect to the minimum sum CRBs of (71) and (72) can be found by searching over M .

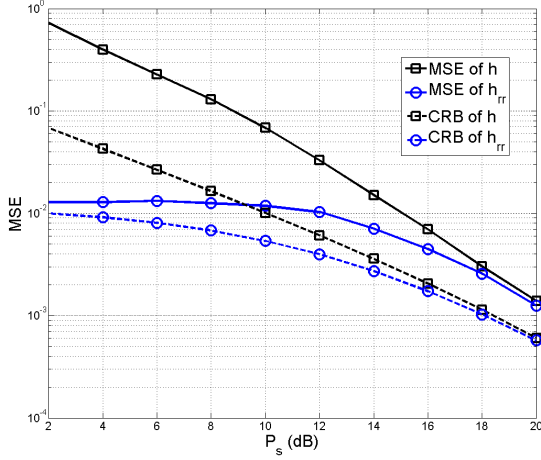


Fig. 1. Performance of ML estimator compared with CRB

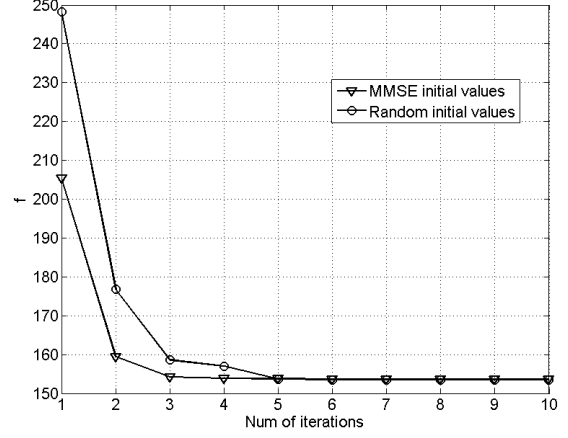


Fig. 2. Number of iterations to convergence for different initial values

VII. NUMERICAL RESULTS

We first simulate the performance of the proposed ML estimator and compare it with the corresponding CRBs. We set $P_r = 30$ dB and $\sigma_{rr}^2 = -10$ dB. For each block we estimate the channels and calculate the mean squared error (MSE) which is averaged over multiple independent realizations of the channels. The training length is $N = 100$.

In Figure 1, we compare the MSEs of h and θ to their CRBs. For h , we obtain the simulated MSEs of h_x and h_y because our optimization only deals with real numbers. To make a fair comparison with its CRB which is derived for complex numbers, we use the fact that the MSE of h is the sum of the MSEs of its real and imaginary parts. The comparisons for θ are similar. The MSE and CRB of h decrease with the relay transmit power P_s getting increased. For θ , the MSE does not decrease when P_s is less than 10 dB. The MSE for θ is also affected by the relay power P_r . It can be seen analytically from (49) that when the amplitude of P_s is close to that of $P_r\sigma_{rr}^2$, the decrease in α is apparent which leads to a decrease in the CRB.

Figure 2 illustrates the convergence speed of the objective function f for different initialization methods, namely, random initialization and MMSE-based initialization. We calculate the average of f in each step for the same h and θ . We observe that with MMSE-based initialization, the objective function converges in 3 iterations while it needs 2 more iterations to converge with random

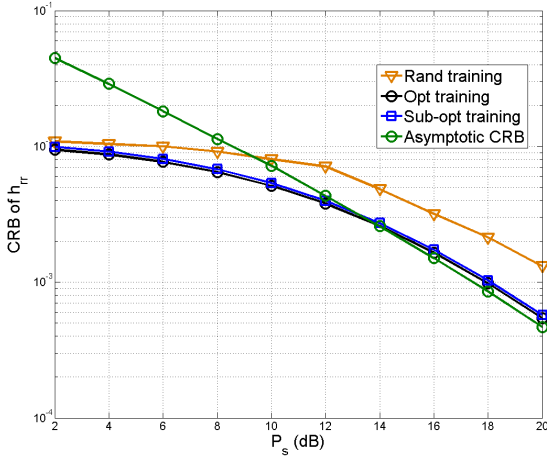


Fig. 3. Comparison of optimal, approximately optimal, and random training sequences

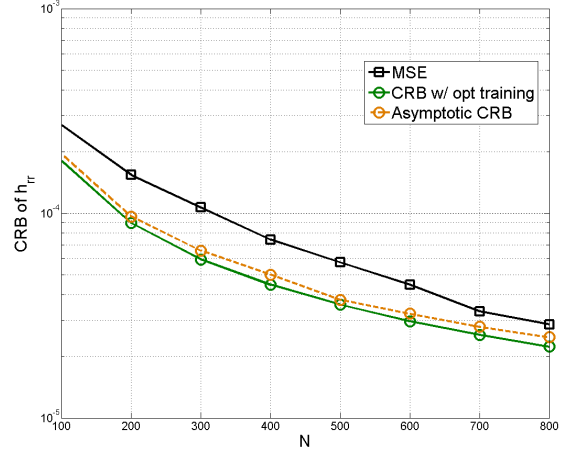


Fig. 4. Effect of training length N to the CRB

initialization. Thus MMSE-based initialization increases the convergence speed of the algorithm.

We compare the CRBs of θ with different training sequences in Figure 3. The CRBs with optimal and approximately optimal training sequence outperform the one with i.i.d Bernoulli symbols. The optimal and approximately optimal curves are almost overlapped. We observe that for different θ values, the roots calculated from the 8th order equation does not fall in the interval $[(1 - |\theta|)^2, (1 + |\theta|)^2]$ discussed in Section IV-B. Thus the optimal solution is on the boundary values which is consistent with the approximately optimal solution. Therefore the simulation results of the optimal and the approximately optimal cases are very close.

Figure 4 shows the influences of training length N on the simulated CRB and the CRB calculated asymptotically, with the optimal training sequence. As N increases, the MSE gradually gets close to the CRB as we expect, since the ML estimator is asymptotically efficient (when N goes to infinity) [26]. The asymptotic CRB which is obtained by (44) is close to the simulated CRB, which shows the efficacy of (44).

In Figure 5, we compare the performance of different detectors with channel tap length $L = 3$. The Viterbi equalizer uses both the CSI for h and θ for equalizing and whitening the noise. The matched filter corresponds to the strongest path of the multi-path channel. The strongest path is related to h and we need θ to whiten the noise. Thus the MF with whitened noise case still uses

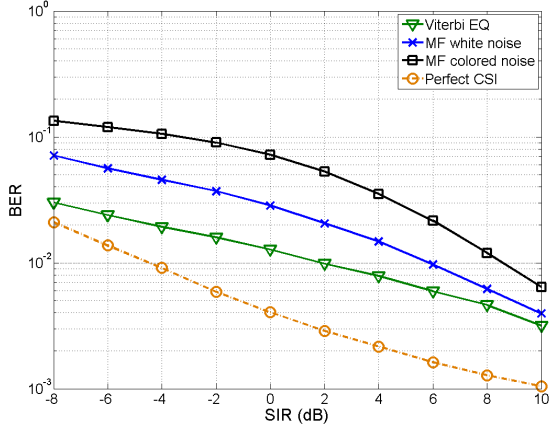
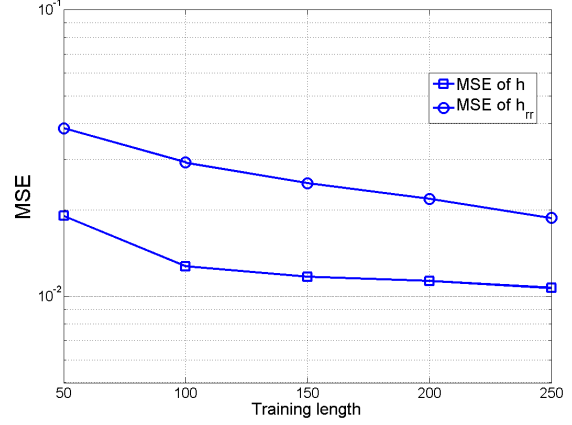


Fig. 5. BER comparison of different detectors

Fig. 6. MSE with increasing N in frequency selective case

both CSI. The MF with colored noise case only uses h . We can see that the Viterbi equalizer cancels the RSI and outperforms the MF, in the white noise case. White noise case of MF is better than the colored noise case. This illustrates the benefits of estimating and canceling the RSI.

Figures 6 and 7 show the MSE of the two extensions. The MSE in these two figures are calculated by comparing the estimates and the exact channels. Note that the estimator are derived by the asymptotic channels \mathbf{H}_f and \mathbf{H}_m which are the approximation of the exact channels for frequency selective case, and multi-relay case respectively. Thus comparing the estimates to the exact channels is a worst-case comparison. Figure 6 shows that in the frequency selective case, the MSE reduces with the training length N increasing, implying that the asymptotic approximation gets closer to the exact channels.

Figure 7 shows the MSEs of z_M and h_{M2} compared with their CRBs in the multi-relay case. Specifically, the total path loss between the source and the relay is $K = -60$ dB and path-loss exponent is $\gamma = 3.71$ for the outdoor environment. As M increases, the MSE of z_M increases because more noise and interference are added. On the other hand, the MSE of h_{M2} decreases, since the RSI channel is easier to estimate as the RSI gets stronger. The asymptotic CRBs derived by (71) and (72) are close to the simulated CRBs. Since M is an integer and not large, one can search over the best M by using (71) and (72).

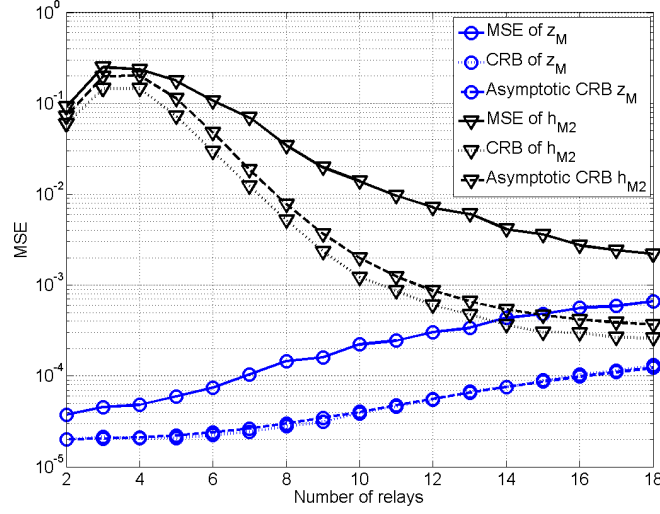


Fig. 7. CRB for multiple relays

VIII. CONCLUSION

We propose an ML channel estimator in FD relays to estimate the end-to-end channel as well as the RSI channel at the destination. The log-likelihood function is maximized through the BFGS algorithm. The algorithm is initialized by a linear MMSE estimator to prevent local minima and increase the convergence speed. The corresponding CRBs are derived to evaluate the accuracy of the estimates. By using Szegő's theorems, we show that the optimal training sequence is a sinusoid. To find the frequency, we minimize the CRBs and propose the corresponding optimal training sequence and a practical approximately optimal training sequence. Extensions of our estimation method to frequency selective and multi-relay case are also considered.

APPENDIX I

DERIVATIVES WITH RESPECT TO NUISANCE PARAMETER d

First we show that $E \left[\frac{\partial f}{\partial d} \right] = 0$. Consider the derivative of f with respect to d ,

$$\frac{\partial f}{\partial d} = -d^* \text{tr} \left(C^{-1} \mathbf{H}_\theta \mathbf{H}_\theta^H \right) + d^* \text{tr} \left((\mathbf{y} - \boldsymbol{\mu})^H C^{-1} \mathbf{H}_\theta \mathbf{H}_\theta^H C^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right). \quad (73)$$

Then take the expectation with respect to noise,

$$\mathbb{E} \left[\frac{\partial f}{\partial d} \right] = -d^* \text{tr} (\mathbf{C}^{-1} \mathbf{H}_\theta \mathbf{H}_\theta^H) + d^* \text{tr} (\mathbb{E}[(\mathbf{y} - \boldsymbol{\mu})^H (\mathbf{y} - \boldsymbol{\mu})] \mathbf{C}^{-1} \mathbf{H}_\theta \mathbf{H}_\theta^H \mathbf{C}^{-1}) \quad (74)$$

$$= -d^* \text{tr} (\mathbf{C}^{-1} \mathbf{H}_\theta \mathbf{H}_\theta^H) + d^* \text{tr} (\mathbf{C} \mathbf{C}^{-1} \mathbf{H}_\theta \mathbf{H}_\theta^H \mathbf{C}^{-1}) = 0. \quad (75)$$

Thus, (14) is proved.

Then we show that (15) holds. We have

$$\frac{\partial^2 f}{\partial h_x \partial d_x} = 2 \text{Re} [(\mathbf{y} - \boldsymbol{\mu})^H \mathbf{C}^{-1} (2d_x \mathbf{H}_\theta \mathbf{H}_\theta^H) \mathbf{C}^{-1} \mathbf{H} \mathbf{x}] \quad (76)$$

$$= 2 \text{Re} \left[(\mathbf{y} - \boldsymbol{\mu})^H 2d_x \mathbf{T}_N \left(\frac{|t(\lambda)|^2}{(|d|^2 |t(\lambda)|^2 + 1)^2} \right) \mathbf{H} \mathbf{x} \right] \quad (77)$$

$$\approx 4d_x \text{Re} \left[(\mathbf{y} - \boldsymbol{\mu})^H \mathbf{T}_N \left(\frac{1}{|d|^4} \right) \mathbf{H} \mathbf{x} \right] \quad (78)$$

$$= \frac{4d_x}{|d|^4} \text{Re} [(\mathbf{y} - \boldsymbol{\mu})^H \mathbf{H} \mathbf{x}]. \quad (79)$$

Equation (77) is obtained by using the asymptotic properties of Toeplitz matrices. We have the approximation in (78) since $|d|^2 \geq |1 - \theta e^{j\lambda}|^2$. Simplify (78) to (79) using $\mathbf{T}_N \left(\frac{1}{|d|^4} \right) = \frac{1}{|d|^4} \mathbf{I}_N$.

Next we need to simplify $\frac{\partial f}{\partial \theta_x}$ by the asymptotic properties of Toeplitz matrices.

$$\begin{aligned} \frac{\partial f}{\partial \theta_x} &= |d|^2 \text{tr} (\mathbf{C}^{-1} (\mathbf{B} \mathbf{H}^H + \mathbf{H} \mathbf{B}^H)) - |d|^2 (\mathbf{y} - \boldsymbol{\mu})^H \mathbf{C}^{-1} (\mathbf{B} \mathbf{H}^H - \mathbf{H} \mathbf{B}^H) \mathbf{C}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \\ &\quad - \text{Re} [(\mathbf{y} - \boldsymbol{\mu})^H \mathbf{C}^{-1} h \mathbf{B} \mathbf{x}] \end{aligned} \quad (80)$$

$$\begin{aligned} &= |d|^2 N \int_0^{2\pi} \frac{1}{|d|^2} \text{Re}[e^{j\lambda} - \theta^*] d\lambda - |d|^2 (\mathbf{y} - \boldsymbol{\mu})^H \mathbf{T}_N \left(\frac{2 \text{Im}[g(\lambda) t^*(\lambda)]}{(|d|^2 |t(\lambda)|^2 + 1)^2} \right) (\mathbf{y} - \boldsymbol{\mu}) \\ &\quad - \text{Re} [(\mathbf{y} - \boldsymbol{\mu})^H \mathbf{C}^{-1} h \mathbf{B} \mathbf{x}] \end{aligned} \quad (81)$$

$$\approx -2\pi N \text{Re}[\theta^2] - \frac{2 \text{Im}[e^{j\lambda} - \theta^*]}{|d|^2} (\mathbf{y} - \boldsymbol{\mu})^H (\mathbf{y} - \boldsymbol{\mu}) - \text{Re} [(\mathbf{y} - \boldsymbol{\mu})^H \mathbf{C}^{-1} h \mathbf{B} \mathbf{x}]. \quad (82)$$

Equation (82) is obtained by approximating the second term in (81). Note that the first term in (82) is not a function of d_x and the third term is similar to (77). Thus we have

$$\frac{\partial^2 f}{\partial \theta_x \partial d_x} = \frac{4d_x \text{Im}[e^{j\lambda} - \theta^*]}{|d|^3} (\mathbf{y} - \boldsymbol{\mu})^H (\mathbf{y} - \boldsymbol{\mu}) - \frac{4d_x}{|d|^4} \text{Re} [(\mathbf{y} - \boldsymbol{\mu})^H h \mathbf{B} \mathbf{x}]. \quad (83)$$

The two second order derivative terms in (79) and (83) both contains factors $\frac{1}{|d|^4}$ and $\frac{1}{|d|^3}$. It can be shown that $d_x \ll |d|^3$, which shows that (79) and (83) are very small and approximately zeros.

APPENDIX II

COEFFICIENTS OF $G'(x)$

To solve $G'(x) = 0$ we only need to solve the numerator of $G'(x)$ equals to zero which is given by

$$\begin{aligned} & |h|^2 \frac{2x^2 + \alpha x}{\theta_R^2} + |h|^2 \frac{x^2 + \alpha x}{8|\theta|^4} \cdot \frac{-\sqrt{\Delta} + m + m|\theta|^2 - mx}{\sqrt{\Delta}} (1 - |\theta|^2 + m\sqrt{\Delta} - x) \\ & - |h|^2 \frac{2\alpha^2 + 3x}{16|\theta|^4} (1 - |\theta|^2 + m\sqrt{\Delta} - x)^2 = 0. \end{aligned} \quad (84)$$

where $m = \theta_I/\theta_R$ and $\Delta = 4|\theta|^2 - (1 + |\theta|^2 - x)^2$. $|h|^2$ can be canceled from (84) and thus the solution is not related to h . Move all the terms containing $\sqrt{\Delta}$ to the left hand side and the others to the right hand side, and take a square of both sides, we have

$$\begin{aligned} & \Delta \left[\frac{x(\alpha + x)(-1 + |\theta|^2 + m^2 + m^2|\theta|^2 + x - m^2x) + m(2\alpha^2 + 3x)(-1 + |\theta|^2 + x)}{|\theta|^4} + \frac{\alpha x + 2x^2}{R^2} \right]^2 \\ & = \left[\frac{m(\alpha x + 2x^2)[(1 - x)^2 - |\theta|^4 - \Delta] - (\alpha^2 + \frac{3}{2}x)[(1 - |\theta|^2 - x)^2 + m\Delta]}{|\theta|^4} \right]^2. \end{aligned} \quad (85)$$

Simplify both sides, we have

$$\begin{aligned} \text{LHS} &= (-x^2 + d_1x + d_2)(e_2x^3 + e_3x^2 + e_4x + e_5)^2, \\ \text{RHS} &= (f_1x^4 + f_2x^3 + f_3x^2 + f_4x + f_5)^2, \end{aligned} \quad (86)$$

where

$$\begin{aligned} d_1 &= 2(1 + |\theta|^2), & d_2 &= -(1 - |\theta|^2)^2, & e_1 &= -1 + |\theta|^2 + m^2 + m^2|\theta|^2, \\ e_2 &= 1 + 3m - m^2, & e_5 &= 2\alpha^2m(|\theta|^2 - 1), \\ e_3 &= \alpha - \alpha m^2 + e_1 - 2\alpha^2m + \frac{2|\theta|^4}{R^2}, & e_4 &= \alpha e_1 + 3m|\theta|^3 - 3m + \frac{\alpha|\theta|^4}{R^2}, \end{aligned}$$

and

$$\begin{aligned}
f_1 &= 2m, & f_2 &= \frac{3}{2} - \frac{3}{2}m^2 - 4m - 2|\theta|^2m + 2\alpha m, \\
f_3 &= 2m - 2|\theta|^2m - 4\alpha m - 2\alpha|\theta|^2m + \alpha^2 - \alpha^2m^2 + 3m^2 + 3m^2|\theta|^2 + 3 + 3|\theta|^2, \\
f_4 &= 2\alpha m - 2\alpha|\theta|^2m + \frac{3}{2}(1 - m^2)(1 - |\theta|^2)^2 + 2\alpha(m^2 + m^2|\theta|^2 - 1 + |\theta|^2), \\
f_5 &= \alpha^2(1 - m^2)(1 - |\theta|^2)^2.
\end{aligned} \tag{87}$$

Finally, the coefficients are

$$\begin{aligned}
x^8 &: f_1^2 + e_2^2, \\
x^7 &: 2f_1f_2 - d_1e_2^2 + 2e_2e_3, \\
x^6 &: (f_2^2 + 2f_1f_3) - (d_1e_2^2 + 2d_1e_2e_3 - e_3^2 - 2e_2e_4), \\
x^5 &: (2f_2f_3 + 2f_1f_4) - (2d_2e_2e_3 + d_1e_3^2 + 2d_1e_2e_4 - 2e_3e_4 - 2e_2e_5), \\
x^4 &: (f_3^2 + 2f_2f_4 + 2f_1f_5) - (d_2e_3^2 + 2d_2e_2e_4 + 2d_1e_3e_4 - e_4^2 + 2e_2e_5 - 2e_3e_5), \\
x^3 &: (2f_3f_4 + 2f_2f_5) - (2d_2e_3e_4 + d_1e_4^2 + 2d_2e_2e_5 + 2d_1e_3e_5 - 2e_4e_5), \\
x^2 &: (f_4^2 + 2f_3f_5) - (d_2e_4^2 + 2d_2e_3e_5 + 2d_2e_4e_5 - e_5^2), \\
x &: 2f_4f_5 - (2d_2e_4e_5 + d_1e_5^2), \\
x^0 &: f_5^2 - d_2e_5^2.
\end{aligned} \tag{88}$$

APPENDIX III

THE GRADIENTS USED IN THE BFGS ALGORITHM

Now we derive the gradients of f with respect to θ_x and θ_y which are used in the BFGS algorithm. The gradients for both real and imaginary parts are needed as inputs of the algorithm. For θ , we first obtain the derivative of \mathbf{H}_θ with respect to θ , denoted as \mathbf{B}_θ , which is also an $N \times N$ Toeplitz matrix with first column $[0, 1, 2\theta, \dots, (L-1)\theta^{L-2}, 0, \dots, 0]^T$ and first row $\mathbf{0}_{N \times 1}^T$.

Both \mathbf{C} and $\boldsymbol{\mu}$ contain θ , therefore there are three terms in its gradient. We have

$$\begin{aligned} \nabla f_{\theta_x} = & \text{tr}(\alpha^2 |d|^2 \sigma_n^2 \mathbf{C}^{-1} (\mathbf{B}_\theta \mathbf{H}_\theta^H + \mathbf{H}_\theta \mathbf{B}_\theta^H)) - 2\text{Re}[(\mathbf{y} - \boldsymbol{\mu})^H \mathbf{C}^{-1} h \mathbf{B}_\theta \mathbf{x}] \\ & - \alpha^2 |d|^2 \sigma_n^2 (\mathbf{y} - \boldsymbol{\mu})^H \mathbf{C}^{-1} (\mathbf{B}_\theta \mathbf{H}_\theta^H + \mathbf{H}_\theta \mathbf{B}_\theta^H) \mathbf{C}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \end{aligned} \quad (89)$$

$$\begin{aligned} \nabla f_{\theta_y} = & \text{tr}(j\alpha^2 |d|^2 \sigma_n^2 \mathbf{C}^{-1} (\mathbf{B}_\theta \mathbf{H}_\theta^H - \mathbf{H}_\theta \mathbf{B}_\theta^H)) - 2\text{Re}[(\mathbf{y} - \boldsymbol{\mu})^H \mathbf{C}^{-1} j h \mathbf{B}_\theta \mathbf{x}] \\ & - j\alpha^2 |d|^2 \sigma_n^2 (\mathbf{y} - \boldsymbol{\mu})^H \mathbf{C}^{-1} (\mathbf{B}_\theta \mathbf{H}_\theta^H - \mathbf{H}_\theta \mathbf{B}_\theta^H) \mathbf{C}^{-1} (\mathbf{y} - \boldsymbol{\mu}). \end{aligned} \quad (90)$$

REFERENCES

- [1] M. Heino, D. Korpi, T. Huusari, E. Antonio-Rodriguez, S. Venkatasubramanian, T. Riihonen, L. Anttila, C. Icheln, K. Haneda, and R. Wichman, "Recent advances in antenna design and interference cancellation algorithms for in-band full duplex relays," *IEEE Communications Magazine*, vol. 53, no. 5, pp. 91–101, 2015.
- [2] S.-K. Hong, J. Brand, J. Choi, M. Jain, J. Mehlman, S. Katti, and P. Levis, "Applications of self-interference cancellation in 5G and beyond," *Communications Magazine, IEEE*, vol. 52, no. 2, pp. 114–121, 2014.
- [3] A. Sabharwal, P. Schniter, D. Guo, D. W. Bliss, S. Rangarajan, and R. Wichman, "In-band full-duplex wireless: Challenges and opportunities," *Selected Areas in Communications, IEEE Journal on*, vol. 32, no. 9, pp. 1637–1652, 2014.
- [4] J. Ma, G. Y. Li, J. Zhang, T. Kuze, and H. Iura, "A new coupling channel estimator for cross-talk cancellation at wireless relay stations," in *Global Telecommunications Conference, 2009. GLOBECOM 2009. IEEE*. IEEE, 2009, pp. 1–6.
- [5] A. Masmoudi and T. Le-Ngoc, "A maximum-likelihood channel estimator for self-interference cancellation in full-duplex systems," *Vehicular Technology, IEEE Transactions on*, 2015, accepted, In press.
- [6] A. Koohian, H. Mehrpouyan, M. Ahmadian, and M. Azarbad, "Bandwidth efficient channel estimation for full duplex communication systems," in *Communications (ICC), 2015 IEEE International Conference on*. IEEE, pp. 4710–4714.
- [7] T. Riihonen, S. Werner, and R. Wichman, "Mitigation of loopback self-interference in full-duplex MIMO relays," *Signal Processing, IEEE Transactions on*, vol. 59, no. 12, pp. 5983–5993, 2011.
- [8] M. Duarte, C. Dick, and A. Sabharwal, "Experiment-driven characterization of full-duplex wireless systems," *Wireless Communications, IEEE Transactions on*, vol. 11, no. 12, pp. 4296–4307, 2012.
- [9] T. Riihonen, S. Werner, and R. Wichman, "Residual self-interference in full-duplex MIMO relays after null-space projection and cancellation," in *Signals, Systems and Computers (ASILOMAR), 2010 Conference Record of the Forty Fourth Asilomar Conference on*. IEEE, 2010, pp. 653–657.
- [10] T. M. Kim and A. Paulraj, "Outage probability of amplify-and-forward cooperation with full duplex relay," in *Wireless Communications and Networking Conference (WCNC), 2012 IEEE*. IEEE, 2012, pp. 75–79.
- [11] L. Jimenez Rodriguez, N. H. Tran, and T. Le-Ngoc, "Performance of full-duplex af relaying in the presence of residual self-interference," *Selected Areas in Communications, IEEE Journal on*, vol. 32, no. 9, pp. 1752–1764, 2014.
- [12] —, "Optimal power allocation and capacity of full-duplex af relaying under residual self-interference," *Wireless Communications Letters, IEEE*, vol. 3, no. 2, pp. 233–236, 2014.

- [13] X. Cheng, B. Yu, X. Cheng, and L. Yang, "Two-way full-duplex amplify-and-forward relaying," in *Military Communications Conference, MILCOM 2013-2013 IEEE*. IEEE, 2013, pp. 1–6.
- [14] F. S. Tabataba, P. Sadeghi, C. Hucher, and M. R. Pakravan, "Impact of channel estimation errors and power allocation on analog network coding and routing in two-way relaying," *Vehicular Technology, IEEE Transactions on*, vol. 61, no. 7, pp. 3223–3239, 2012.
- [15] D. Kim, H. Ju, S. Park, and D. Hong, "Effects of channel estimation error on full-duplex two-way networks," *Vehicular Technology, IEEE Transactions on*, vol. 62, no. 9, pp. 4666–4672, 2013.
- [16] G. Zheng, "Joint beamforming optimization and power control for full-duplex MIMO two-way relay channel," *Signal Processing, IEEE Transactions on*, vol. 63, no. 3, pp. 555–566, 2015.
- [17] X. Li, C. Tepedelenlioğlu, and H. Şenol, "Channel estimation for residual self-interference in full duplex amplify-and-forward two-way relays," *Wireless Communications, IEEE Transactions on*, vol. 16, no. 8, pp. 4970–4983, 2017.
- [18] X. Xiong, X. Wang, T. Riihonen, and X. You, "Channel estimation for full-duplex relay systems with large-scale antenna arrays," *Wireless Communications, IEEE Transactions on*, vol. 15, no. 10, pp. 6925–6938, 2016.
- [19] X. Li and C. Tepedelenlioğlu, "Maximum likelihood channel estimation for residual self-interference cancellation in full duplex relay," in *Signals, Systems and Computers (ASILOMAR), 2015 Conference Record of the Forty Ninth Asilomar Conference on*. IEEE, 2015, pp. 807–811.
- [20] H. Q. Ngo, H. A. Suraweera, M. Matthaiou, and E. G. Larsson, "Multipair full-duplex relaying with massive arrays and linear processing," *IEEE Journal on Selected Areas in Communications*, vol. 32, no. 9, pp. 1721–1737, 2014.
- [21] T. Riihonen, S. Werner, and R. Wichman, "Optimized gain control for single-frequency relaying with loop interference," *Wireless Communications, IEEE Transactions on*, vol. 8, no. 6, pp. 2801–2806, 2009.
- [22] K. P. Murphy, *Machine Learning: A Probabilistic Perspective*. MIT press, 2012.
- [23] U. Grenander and G. Szegő, *Toeplitz Forms and Their Applications*. Univ of California Press, 2001, vol. 321.
- [24] R. M. Gray, *Toeplitz and Circulant Matrices: A Review*. Now Publishers Inc, 2006.
- [25] A. Goldsmith, *Wireless communications*. Cambridge university press, 2005.
- [26] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*. Prentice Hall, 1993.