This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

arXiv:1710.10830v1 [cs.IT] 30 Oct 2017

# A Framework for Over-the-air Reciprocity Calibration for TDD Massive MIMO Systems

Xiwen JIANG, Alexis Decurninge, Kalyana Gopala, Florian Kaltenberger, *Member, IEEE,*
Maxime Guillaud, *Senior Member, IEEE,* Dirk Slock, *Fellow, IEEE,* and Luc Deneire, *Member, IEEE*

*Abstract*—One of the biggest challenges in operating massive multiple-input multiple-output systems is the acquisition of accurate channel state information at the transmitter. To take up this challenge, time division duplex is more favorable thanks to its channel reciprocity between downlink and uplink. However, while the propagation channel over the air is reciprocal, the radio-frequency front-ends in the transceivers are not. Therefore, calibration is required to compensate the RF hardware asymmetry.

Although various over-the-air calibration methods exist to address the above problem, this paper offers a unified representation of these algorithms, providing a higher level view on the calibration problem, and introduces innovations on calibration methods. We present a novel family of calibration methods, based on antenna grouping, which improve accuracy and speed up the calibration process compared to existing methods. We then provide the Cramér-Rao bound as the performance evaluation benchmark and compare maximum likelihood and least squares estimators. We also differentiate between coherent and non-coherent accumulation of calibration measurements, and point out that enabling non-coherent accumulation allows the training to be spread in time, minimizing impact to the data service. Overall, these results have special value in allowing to design reciprocity calibration techniques that are both accurate and resource-effective.

*Index Terms*—Massive MIMO, TDD, channel reciprocity calibration.

## I. INTRODUCTION

Massive multiple-input multiple-output (MIMO) is a promising air interface technology for the next generation of wireless communications. With large number of antennas installed at the base station (BS) simultaneously serving multiple user equipments (UEs), massive MIMO can dramatically improve the spectral efficiency of cellular networks [1], [2].

For downlink (DL), one of the fundamental challenges to fully realize the potential of massive MIMO is the acquisition of accurate channel state information at the transmitter (CSIT). Time division duplex (TDD) thus attracts great attention

from the research community as it enjoys channel reciprocity between DL and uplink (UL), thanks to which the BS can obtain the CSIT from the channel estimation in the UL. In fact, traditional ways to get CSIT from UE feedback becomes infeasible when the antenna array size at the BS scales up, because of the heavy signaling overhead it incurs in the UL.

Channel reciprocity in TDD systems refers to the fact that the physical over-the-air (OTA) channels are the same for UL and DL [3], [4] within channel coherence time. However, the channel as seen by the digital baseband processor contains not only the physical OTA channel but also radio frequency (RF) front-ends, including the hardware from digital-to-analog converter (DAC) to transmit antennas at the transmitter (Tx) and the corresponding part, from receiving antennas to analog-to-digital converter (ADC), at the receiver (Rx). The various impairments to reciprocity can be due to manufacturing variability in the power amplifiers and low-noise amplifiers, different cable lengths across the antennas, imperfect clock synchronization, duplexer response, etc. Due to these, the hardware in the Tx and Rx RF chains are, in general, not identical, and therefore the channel from a digital signal processing point of view is not reciprocal. If not taken into account, these hardware-related asymmetries will cause inaccuracy in the CSIT estimation and, as a consequence, seriously degrade the DL beamforming performance [5]–[8].

In order to compensate the hardware asymmetry and restore channel reciprocity, calibration techniques are needed. This topic has been explored long before the advent of massive MIMO. In [9]–[13], it is suggested to add additional hardware components in transceivers which are dedicated to calibration. This method (which we refer to as *absolute* calibration) consists in compensating the Tx and Rx RF asymmetry independently in each transceiver; however this does not appear to be a cost-effective solution. [14]–[17] thus put forward "relative" calibration schemes[1], where the calibration coefficients are estimated using signal processing methods based on OTA bi-directional channel estimation between BS and UE. Since hardware properties can be expected to evolve slowly, and these coefficients can be obtained in the initialization phase of the system (calibration phase), they can be used later together with the instantaneous UL channel estimate to obtain downlink CSIT.

With the advent of massive MIMO, traditional relative calibration methods are challenged, because they require the

X. Jiang, K. Gopala, F. Kaltenberger and D. Slock are with Communication Systems Department, EURECOM. (e-mail: {xiwen.jiang, kalyana.gopala, florian.kaltenberger, dirk.slock}@eurecom.fr)

A. Decurninge and M. Guillaud are with the Mathematical and Algorithmic Sciences Lab, Paris Research Center, Huawei Technologies France. (e-mail: {alexis.decurninge, maxime.guillaud}@huawei.com)

L. Deneire is with Laboratoire I3S, Université de Nice Sophia Antipolis. (e-mail: luc.deneire@unice.fr)

---

[1]The term *relative* indicates here that the calibration coefficients relate the UL and DL digital channels, as opposed to absolute calibration which relates digital domain and propagation domain versions of a channel.

UE to feed back a large amount of DL CSI for all BS antennas. It was observed in [18] that the calibration factor at the BS side is the same for all channels from the BS to any UE. This was exploited in [18] to determine the BS side calibration factor of a secondary BS with the cooperation of a secondary UE, allowing beamforming with zero-forcing to a primary UE without its collaboration. This idea was then pushed further in a number of OTA self-calibration approaches which only require the exchange of OTA training signals between elements of the BS array. Indeed, for optimizing multi-user massive MIMO systems, the asymmetry in the number of antennas between the BS and the UEs means that most of the massive MIMO multi-user multiplexing gain can be achieved through BS-side only calibration [19], [20]. These OTA self-calibration approaches have the advantage that, unlike classical single-link relative calibration, no CSI feedback is involved, since all the elements of the BS array are already connected to the same baseband signal processor. Such "single-side" or "internal" calibration methods were proposed in [21]–[25]. In [21], the authors reported on the massive MIMO Argos prototype, where calibration is performed OTA with the help of a reference antenna. By performing bi-directional transmission between the reference antenna and the rest of the antenna array, it is possible to estimate the calibration coefficients up to a common scalar ambiguity which will not influence the final DL beamforming capability. The Argos calibration approach however is sensitive to the location of the reference antenna, and as one of the consequences, is not suitable for distributed massive MIMO. This concern motivated the introduction of a method (Rogalin *et al.* in [22]) whereby calibration is not performed w.r.t. a reference antenna[2]. It has the spirit of distributed algorithms, making it a good calibration method for antenna arrays having a distributed topology. Note that it can also be applied to colocated massive MIMO, as in the LuMaMi massive MIMO prototype [26] where a weighted version of the estimator presented in [23] is used, whereas a Maximum Likelihood (ML) estimator is presented in [24]. Moreover, a fast calibration method named Avalanche was proposed in [25]; its principle is to use a calibrated sub-array to calibrate uncalibrated elements. The calibrated array thus grows during the calibration process in a way similar to the avalanche phenomenon.

Among other relevant works, we refer to [27]–[31]. In [27], the author provides an idea to perform system health monitoring on the calibrated reciprocity. Under the assumption that the majority of calibration coefficients stay calibrated and only a minority of them change, the authors propose a compressed sensing enabled detection algorithm to find out which calibration coefficient has changed based on the sparsity in the vector representing the coefficient change. In [28], a calibration method dedicated to maximum ratio transmission (MRT) is proposed. Experimental data about the calibration coefficients are reported in [21], [29]–[31], giving an insight on how the impairments evolve in the time and frequency domains

as well as with the temperature, and about the hardware properties behind this effect.

In the present article, we introduce a unified framework to represent different existing calibration methods. Although they appear at first sight to be different, we reveal that all existing calibration methods can be modeled under a general pilot based calibration framework; different ways to partition the array into transmit and receive elements during successive training phases yield different schemes. The unified representation shows the relationship between these methods and provides alternative ways to obtain corresponding estimators. As this framework gives a general and high level understanding of the TDD calibration problem in massive MIMO systems, it opens up possibilities for new calibration methods. As an example, we present a novel family of calibration schemes based on antenna grouping, which can greatly speed up the calibration process with respect to the classical approaches. We will show that our proposed method greatly outperforms the Avalanche method [25] in terms of calibration accuracy, yet it is equally fast. In order to evaluate the performance of calibration schemes, we derive the Cramér-Rao bounds (CRB) of the accuracy of calibration coefficients estimation. Another important contribution of this work is the introduction of non-coherent accumulation of the measurements used for calibration. We will see that calibration does not necessarily have to be performed in an intensive manner during a single channel coherence interval, but can rather be executed using time resources distributed over a relatively long period. This enables TDD reciprocity calibration to be interleaved with the normal data transmission or reception, leaving it almost invisible for the whole system.

The rest of this paper is organized as follows. Section II describes the basic principles of reciprocity calibration in a TDD based MIMO system. Section III presents the TDD reciprocity system model and introduces our unified framework. Section IV presents how Argos, Rogalin and Avalanche calibration algorithms fit into this model as well as how we can obtain the corresponding estimators. In Section V, we present the fast calibration scheme based on antenna grouping and discuss the minimum number of channel uses it requires to estimate all calibration coefficients. In Section VI, we address the optimal estimation problem of reciprocity calibration parameters, we derive the CRB, propose a maximum likelihood (ML) estimator and compare it with the LS estimator. Section VII is dedicated to non-coherent accumulation of measurements. In Section VIII, we illustrate the performance of the group-based fast calibration method and compare its performance with other calibration algorithms using CRB as the benchmark. Conclusions are drawn in Section IX.

The notation adopted in this paper conforms to the following conventions. Vectors and matrices are denoted in lowercase bold and uppercase bold respectively: $\mathbf{a}$ and $\mathbf{A}$. $(\cdot)^*$, $(\cdot)^T$, $(\cdot)^H$, $(\cdot)^\dagger$ denote element-wise complex conjugate, transpose, Hermitian transpose and Moore-Penrose pseudo inverse, respectively. $\otimes$ and $*$ denotes the Kronecker product operator and the Khatri–Rao product [32], respectively. $\lceil \cdot \rceil$ is the ceiling operator, which rounds a number to the next integer. $\mathrm{diag}\{a_1, a_2, \ldots, a_M\}$ denotes a diagonal matrix with

---

[2]The method in [22] is denoted as "least-squares (LS) calibration", however we will not use this terminology since most calibration techniques proposed in the literature ultimately rely on LS estimation

its diagonal composed of $a_1, a_2, \ldots, a_M$, whereas $\text{vec}(\mathbf{A})$ denotes the vectorization of the matrix $\mathbf{A}$. $\mathbb{C}$ denotes the set of complex numbers.

## II. OTA RECIPROCITY CALIBRATION

In this section, we describe the basic idea of reciprocity calibration in a practical TDD system. Let us consider a system as in Fig. 1, where A represents a BS and B represents a UE, each containing $M_A$ and $M_B$ antennas, respectively. The DL and UL channels flat-fading model (as typically obtained by considering a single subcarrier of a multicarrier system) seen in the digital domain are noted by $\mathbf{H}_{A\to B}$ and $\mathbf{H}_{B\to A}$. Since they are formed by the cascade of the Tx impairments, OTA propagation, and Rx impairments, they can be represented by

$$\begin{cases} \mathbf{H}_{A\to B} = \mathbf{R}_B \mathbf{C}_{A\to B} \mathbf{T}_A, \\ \mathbf{H}_{B\to A} = \mathbf{R}_A \mathbf{C}_{B\to A} \mathbf{T}_B, \end{cases} \quad (1)$$

where matrices $\mathbf{T}_A$, $\mathbf{R}_A$, $\mathbf{T}_B$, $\mathbf{R}_B$ model the response of the transmit and receive RF front-ends, while $\mathbf{C}_{A\to B}$ and $\mathbf{C}_{B\to A}$ model the OTA propagation channels, respectively from A to B and from B to A. The dimension of $\mathbf{T}_A$ and $\mathbf{R}_A$ are $M_A \times M_A$, whereas that of $\mathbf{T}_B$ and $\mathbf{R}_B$ are $M_B \times M_B$. The diagonal elements in these matrices represent the linear effects attributable to the impairments in the transmitter and receiver parts of the RF front-ends respectively, whereas the off-diagonal elements correspond to RF crosstalk and antenna mutual coupling[3]. It is worth noting that although transmitting and receiving antenna mutual coupling is not generally recip-rocal [34], theoretical modeling [12] and experimental results [21], [24], [31] both show that in practice, RF crosstalk and antenna mutual coupling can be ignored for the purpose of reciprocity calibration, which implies that $\mathbf{T}_A$, $\mathbf{R}_A$, $\mathbf{T}_B$, $\mathbf{R}_B$ can safely be assumed to be diagonal.

Assuming the system is operating in TDD mode, the OTA channel responses enjoy reciprocity within the channel coherence time, i.e., $\mathbf{C}_{A\to B} = \mathbf{C}_{B\to A}^T$. Therefore, we obtain the following relationship between the channels measured in both directions:

$$\begin{aligned} \mathbf{H}_{A\to B} &= \mathbf{R}_B (\mathbf{R}_A^{-1} \mathbf{H}_{B\to A} \mathbf{T}_B^{-1})^T \mathbf{T}_A \\ &= \underbrace{\mathbf{R}_B \mathbf{T}_B^{-T}}_{\mathbf{F}_B^{-T}} \mathbf{H}_{B\to A}^T \underbrace{\mathbf{R}_A^{-T} \mathbf{T}_A}_{\mathbf{F}_A} \\ &= \mathbf{F}_B^{-T} \mathbf{H}_{B\to A}^T \mathbf{F}_A. \end{aligned} \quad (2)$$

A system utilizing OTA reciprocity calibration normally has two phases for its function. Firstly, during the initialization of the system, the calibration process is performed, which consists in estimating $\mathbf{F}_A$ and $\mathbf{F}_B$. Then during the data transmission phase, they are used together with instantaneous measured UL channel $\hat{\mathbf{H}}_{B\to A}$ to estimate $\mathbf{H}_{A\to B}$ according to (2), based on which advanced beamforming algorithms can be performed. Since the calibration coefficients typically

[3]Here, "antenna mutual coupling" is used to describe parasitic effects that two nearby antennas have on each other, when they are either both transmitting or receiving [12], [33]. However, this is different to the channel between transmitting and receiving elements of the same array, which we call the intra-array channel. Note that this differs from the terminology in [24] and [28] where the term mutual coupling is used to denote the intra-array channel.
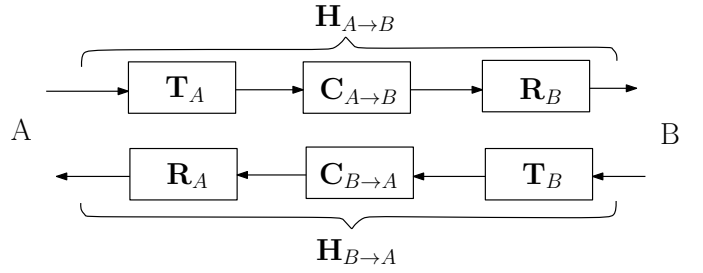


Fig. 1. Reciprocity Model

remain stable [21], the calibration process does not have to be performed very frequently.

Note that the studies in [19], [20] pointed out that in a practical multi-user MIMO system, it is mainly the calibration at the BS side which restores the hardware asymmetry and helps to achieve the multi-user MIMO performance, whereas the benefit brought by the calibration on the UE side is not necessarily justified. We thus, in the sequel, focus on the estimation of $\mathbf{F}_A$, although the framework discussed in the following section is not limited to this case.

## III. GENERAL OTA CALIBRATION FRAMEWORK

### A. Overview and signalling

In this section, we present a general framework for OTA pilot-based reciprocity calibration. Let us consider an antenna array of $M$ elements partitioned into $G$ groups denoted by $A_1, A_2, \ldots, A_G$, as in Fig. 2. Group $A_i$ contains $M_i$ antennas such that $\sum_{i=1}^G M_i = M$. Each group $A_i$ transmits a sequence of $L_i$ pilot symbols, defined by matrix $\mathbf{P}_i \in \mathbb{C}^{M_i \times L_i}$ where the rows correspond to antennas and the columns to successive channel uses. Note that a channel use can be understood as a time slot or a subcarrier in an OFDM-based system, as long as the calibration parameter can be assumed constant over all channel uses. When an antenna group $i$ transmits, all other groups are considered in receiving mode. After all $G$ groups have transmitted, the received signal for each resource block of bidirectional transmission between antenna groups $i$ and $j$ is given by

$$\begin{cases} \mathbf{Y}_{i\to j} = \mathbf{R}_j \mathbf{C}_{i\to j} \mathbf{T}_i \mathbf{P}_i + \mathbf{N}_{i\to j}, \\ \mathbf{Y}_{j\to i} = \mathbf{R}_i \mathbf{C}_{j\to i} \mathbf{T}_j \mathbf{P}_j + \mathbf{N}_{j\to i}, \end{cases} \quad (3)$$

where $\mathbf{Y}_{i\to j} \in \mathbb{C}^{M_j \times L_i}$ and $\mathbf{Y}_{j\to i} \in \mathbb{C}^{M_i \times L_j}$ are received signal matrices at antenna groups $j$ and $i$ respectively when the other group is transmitting. $\mathbf{N}_{i\to j}$ and $\mathbf{N}_{j\to i}$ represent the corresponding received noise matrix. $\mathbf{T}_i$, $\mathbf{R}_i \in \mathbb{C}^{M_i \times M_i}$ and $\mathbf{T}_j$, $\mathbf{R}_j \in \mathbb{C}^{M_j \times M_j}$ represent the effect of the transmit and receive RF front-ends of antenna elements in groups $i$ and $j$ respectively.

The reciprocity property induces that $\mathbf{C}_{i\to j} = \mathbf{C}_{j\to i}^T$, thus for two different groups $1 \le i \ne j \le G$, by eliminating $\mathbf{C}_{i\to j}$ in (3) we have

$$\mathbf{P}_i^T \mathbf{F}_i^T \mathbf{Y}_{j\to i} - \mathbf{Y}_{i\to j}^T \mathbf{F}_j \mathbf{P}_j = \widetilde{\mathbf{N}}_{ij}, \quad (4)$$

where the noise component $\widetilde{\mathbf{N}}_{ij} = \mathbf{P}_i^T \mathbf{F}_i^T \mathbf{N}_{j\to i} - \mathbf{N}_{i\to j}^T \mathbf{F}_j \mathbf{P}_j$, while $\mathbf{F}_i = \mathbf{R}_i^{-T} \mathbf{T}_i$ and $\mathbf{F}_j = \mathbf{R}_j^{-T} \mathbf{T}_j$ are the calibration matrices for groups $i$ and $j$. The calibration matrix $\mathbf{F}$ is diagonal,
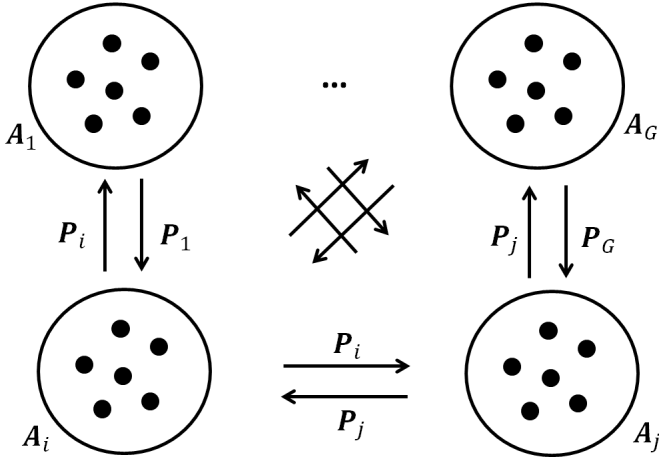
Fig. 2. Bi-directional transmission between antenna groups.

and thus takes the form of $\mathbf{F} = \mathrm{diag}\{\mathbf{F}_1, \mathbf{F}_2, \ldots, \mathbf{F}_G\}$. Note that estimating $\mathbf{F}_i$ or $\mathbf{F}_j$ from (4) for a given pair $(i, j)$ does not exploit all relevant received data. An optimal estimation jointly considering all received signals for all $(i, j)$ will be proposed in Section VI. Note that the proposed framework also allows to consider using only subsets of the received data which corresponds to some of the methods found in the literature.

Let us use $\mathbf{f}_i$ and $\mathbf{f}$ to denote the vectors of the diagonal coefficients of $\mathbf{F}_i$ and $\mathbf{F}$ respectively, i.e., $\mathbf{F}_i = \mathrm{diag}\{\mathbf{f}_i\}$ and $\mathbf{F} = \mathrm{diag}\{\mathbf{f}\}$. This allows us to vectorize (4) into

$$(\mathbf{Y}_{j \to i}^T * \mathbf{P}_i^T)\mathbf{f}_i - (\mathbf{P}_j^T * \mathbf{Y}_{i \to j}^T)\mathbf{f}_j = \widetilde{\mathbf{n}}_{ij}, \tag{5}$$

where $*$ denotes the Khatri–Rao product (or column-wise Kronecker product[4]), where we have used the equality $\mathrm{vec}(\mathbf{A}\,\mathrm{diag}(\mathbf{x})\,\mathbf{B}) = (\mathbf{B}^T * \mathbf{A})\,\mathbf{x}$. Note that, if we do not suppose that every $\mathbf{F}_i$ is diagonal, (5) holds more generally by replacing the Katri–Rao products by Kronecker products and $\mathbf{f}_i$ by $\mathrm{vec}(\mathbf{F}_i)$. Finally, stacking equations (5) for all $1 \leq i < j \leq G$ yields

$$\boldsymbol{\mathcal{Y}}(\mathbf{P})\mathbf{f} = \widetilde{\mathbf{n}}, \tag{6}$$

with $\boldsymbol{\mathcal{Y}}(\mathbf{P})$ defined as

$$\underbrace{\begin{bmatrix} (\mathbf{Y}_{2\to1}^T * \mathbf{P}_1^T) & -(\mathbf{P}_2^T * \mathbf{Y}_{1\to2}^T) & 0 & \cdots \\ (\mathbf{Y}_{3\to1}^T * \mathbf{P}_1^T) & 0 & -(\mathbf{P}_3^T * \mathbf{Y}_{1\to3}^T) & \cdots \\ 0 & (\mathbf{Y}_{3\to2}^T * \mathbf{P}_2^T) & -(\mathbf{P}_3^T * \mathbf{Y}_{2\to3}^T) & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}}_{(\sum_{j=2}^{G} \sum_{i=1}^{j-1} L_i L_j) \times M}.$$

(7)

It is worth noting that this framework is not limited to represent single-side calibration. For UE-aided (relative) calibration, it suffices to set 2 groups such as $A_1$ and $A_2$, representing the BS and the UE, respectively in order to get a full calibration scheme.

[4]With matrices $\mathbf{A}$ and $\mathbf{B}$ partitioned into columns, $\mathbf{A} = [\mathbf{a}_1 \;\; \mathbf{a}_2 \;\; \ldots \;\; \mathbf{a}_M]$ and $\mathbf{B} = [\mathbf{b}_1 \;\; \mathbf{b}_2 \;\; \ldots \;\; \mathbf{b}_M]$ where $\mathbf{a}_i$ and $\mathbf{b}_i$ are column vectors for $i \in 1 \ldots M$, then, $\mathbf{A} * \mathbf{B} = [\mathbf{a}_1 \otimes \mathbf{b}_1 \;\; \mathbf{a}_2 \otimes \mathbf{b}_2 \;\; \ldots \;\; \mathbf{a}_M \otimes \mathbf{b}_M]$ [32].

### B. Parameter identifiability and pilot design

Before proposing an estimator for $\mathbf{f}$, we raise the question of the problem identifiability which corresponds to the fact that (6) admits a unique solution in the noiseless scenario

$$\boldsymbol{\mathcal{Y}}(\mathbf{P})\mathbf{f} = \mathbf{0}. \tag{8}$$

The solution of (8) is defined up to a complex scalar factor $\alpha$, since if $\mathbf{f}$ is a solution, then $\alpha\mathbf{f}$ is also a solution of (8). This indeterminacy can be resolved by fixing one of the calibration parameters, say $f_1 = \mathbf{e}_1^H\mathbf{f} = [1\,0 \cdots 0]\mathbf{f} = 1$ or by a norm constraint, for example $\|\mathbf{f}\| = 1$. Then, the identifiability is related to the dimension of the kernel of $\boldsymbol{\mathcal{Y}}(\mathbf{P})$ in the sense that the problem is fully determined if and only if the kernel of $\boldsymbol{\mathcal{Y}}(\mathbf{P})$ is of dimension 1. Since the true $\mathbf{f}$ is a solution to (8), we know that the rank of $\boldsymbol{\mathcal{Y}}(\mathbf{P})$ is at most $M - 1$. We will assume furthermore in the following that the pilot design is such that the rows of $\boldsymbol{\mathcal{Y}}(\mathbf{P})$ are linearly independent as long as the number of rows is less than $M - 1$. Note that this condition depends on the internal channel realization $\mathbf{C}_{i \to j}$ and on the pilot matrices $\mathbf{P}_i$. However, sufficient conditions of identifiability expressed on these matrices are out of the scope of this paper. Under rows independence, (6) may be read as the following sequence of events:

1) Group 1 broadcasts its pilots to all other groups using $L_1$ channel uses;
2) After group 2 transmits its pilots, we can formulate $L_2 L_1$ equations of the form (5);
3) After group 3 transmits its pilots, we can formulate $L_3 L_1 + L_3 L_2$ equations;
4) After group j transmits its pilots, we can formulate $\sum_{i=1}^{j-1} L_j L_i$ equations.

This process continues until group $G$ finishes its transmission, and the whole calibration process finishes. During this process of transmission by the $G$ antenna groups, we can start forming equations as indicated, that can be solved recursively for subsets of unknown calibration parameters, or we can wait until all equations are formed to solve the problem jointly. By independence of the rows, we can state that the problem is fully determined if and only if

$$\sum_{1 \leq i < j \leq G} L_j L_i \geq M - 1 \,. \tag{9}$$

### C. LS calibration parameter estimation

A typical way to estimate $\mathbf{f}$ consists in solving a LS problem such as

$$\begin{aligned} \hat{\mathbf{f}} &= \arg\min_{\mathbf{f}} \|\boldsymbol{\mathcal{Y}}(\mathbf{P})\,\mathbf{f}\|^2 \\ &= \arg\min_{\mathbf{f}} \sum_{i<j} \|(\mathbf{Y}_{j\to i}^T * \mathbf{P}_i^T)\mathbf{f}_i - (\mathbf{P}_j^T * \mathbf{Y}_{i\to j}^T)\mathbf{f}_j\|^2 \,, \end{aligned} \tag{10}$$

where $\boldsymbol{\mathcal{Y}}(\mathbf{P})$ is defined in (7). This needs to be augmented with a constraint

$$\mathcal{C}(\hat{\mathbf{f}}, \mathbf{f}) = 0, \tag{11}$$

in order to exclude the trivial solution $\hat{\mathbf{f}} = \mathbf{0}$ in (10). The constraint on $\hat{\mathbf{f}}$ may depend on the true parameters $\mathbf{f}$. As we shall see further this constraint needs to be complex valued

(which represents two real constraints). Typical choices for the constraint are

1) Norm plus phase constraint (NPC):

$$\text{norm: } \mathrm{Re}\{\mathcal{C}(\hat{\mathbf{f}}, \mathbf{f})\} = ||\hat{\mathbf{f}}||^2 - c, \ c = ||\mathbf{f}||^2, \quad (12)$$

$$\text{phase: } \mathrm{Im}\{\mathcal{C}(\hat{\mathbf{f}}, \mathbf{f})\} = \mathrm{Im}\{\hat{\mathbf{f}}^H \mathbf{f}\} = 0. \quad (13)$$

2) Linear constraint:

$$\mathcal{C}(\hat{\mathbf{f}}, \mathbf{f}) = \hat{\mathbf{f}}^H \mathbf{g} - c = 0 \ . \quad (14)$$

If we choose the vector $\mathbf{g} = \mathbf{f}$ and $c = ||\mathbf{f}||^2$, then the $\mathrm{Im}\{.\}$ part of (14) corresponds to (13). The most popular linear constraint is the First Coefficient Constraint (FCC), which is (14) with $\mathbf{g} = \mathbf{e}_1$, $c = 1$. The solution of (10), (14) is given by

$$
\begin{aligned}
\hat{\mathbf{f}} &= \arg \min_{\mathbf{f}:\mathbf{f}^H \mathbf{g} = c} ||\boldsymbol{\mathcal{Y}}(\mathbf{P})\,\mathbf{f}||^2 \\
&= \frac{c}{\mathbf{g}^H (\boldsymbol{\mathcal{Y}}(\mathbf{P})^H \boldsymbol{\mathcal{Y}}(\mathbf{P}))^{-1} \mathbf{g}} (\boldsymbol{\mathcal{Y}}(\mathbf{P})^H \boldsymbol{\mathcal{Y}}(\mathbf{P}))^{-1} \mathbf{g} \ .
\end{aligned}
\quad (15)
$$

Assuming a unit norm constraint ((12) with $c = 1$) on the other hand yields

$$\hat{\mathbf{f}}' = \arg \min_{\mathbf{f}: ||\mathbf{f}|| = 1} ||\boldsymbol{\mathcal{Y}}(\mathbf{P})\,\mathbf{f}||^2 = V_{min}(\boldsymbol{\mathcal{Y}}(\mathbf{P})^H \boldsymbol{\mathcal{Y}}(\mathbf{P})), \quad (16)$$

where $V_{min}(\mathbf{X})$ denotes the eigenvector of matrix $\mathbf{X}$ corresponding to its eigenvalue with the smallest magnitude. Then the NPC solution of (10), (12), (13) is $\hat{\mathbf{f}} = \sqrt{c}\, e^{j\phi}\hat{\mathbf{f}}'$ in which the phase $\phi$ is adjusted to satisfy (13), i.e. $\phi = \arg(\hat{\mathbf{f}}'^H \mathbf{f})$ where for any complex number $z = |z|e^{j\arg(z)}$.

## IV. EXISTING CALIBRATION TECHNIQUES

Different choices for the partitioning of the $M$ antennas and the pilots matrices exposed in Section III lead to different calibration algorithms. We will now see how different estimators of the calibration matrix can be derived from (5). In order to ease the description, we assume that the channel is constant during the whole calibration process, this assumption will later be relaxed and discussed in Section VII.

### A. Argos

The Argos calibration method [21] consists in performing bi-directional transmission between a carefully chosen reference antenna and the rest of the antenna array. This can be recast in our framework by considering $G = 2$ sets of antennas, with set $A_1$ containing only the reference antenna ($M_1 = 1$), and set $A_2$ containing all the other antenna elements ($M_2 = M - 1$), as shown in Fig. 3. Firstly, pilot 1 is broadcasted from the reference antenna to all antennas in set $A_2$, thus $L_1 = 1$, $\mathbf{P}_1 = 1$ and $\mathbf{f}_2 = \begin{bmatrix} f_2, \ldots, f_M \end{bmatrix}^T$. Then, antennas in set $A_2$ successively transmit pilot 1 to the reference antenna, thus $L_2 = M - 1$ and $\mathbf{P}_2 = \mathbf{I}_{M-1}$. (5) thus becomes

$$f_1 \mathbf{y}_1 = \mathrm{diag}(\mathbf{y}_2)\mathbf{f}_2 + \tilde{\mathbf{n}}, \quad (17)$$

where $\mathbf{y}_1 = \begin{bmatrix} y_{2 \to 1} & y_{3 \to 1} & \cdots & y_{M \to 1} \end{bmatrix}^T$ and $\mathbf{y}_2 = \begin{bmatrix} y_{1 \to 2} & y_{1 \to 3} & \cdots & y_{1 \to M} \end{bmatrix}^T$ with $y_{i \to j}$ representing the signal transmitted from antenna $i$ and received at antenna $j$. (17) can be decomposed into $M - 1$ independent equations
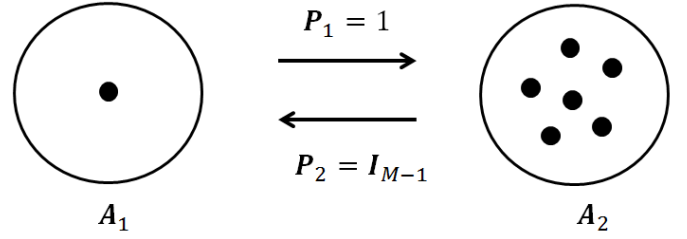


Fig. 3. Argos calibration

as $f_1 y_{i \to 1} = f_i y_{1 \to i} + \tilde{n}_i$, where $\tilde{n}_i$ is the $i^{th}$ element in the noise vector $\tilde{\mathbf{n}}$. The LS estimator for each element is thus

$$f_i = f_1 \frac{y_{i \to 1}}{y_{1 \to i}}, \quad \text{where } i = 2, 3, \ldots, M. \quad (18)$$

### B. Methods based on successive single-antenna transmissions followed by joint estimation

The method from Rogalin et al. presented in [22], [35] and further analyzed in [24] is based on single-antenna transmission at each channel use; all received signals are subsequently taken into account through joint estimation of the calibration parameters. In order to represent this method within the unified framework, we define each set $A_i$ as containing only antenna $i$, i.e., $M_i = 1$ for $1 \le i \le M$, as in Fig. 4. Since we assume that the channel is constant, this calibration procedure can be performed in a way that antennas can broadcast pilot 1 in a round-robin manner to all other antennas. In total, $M$ channel uses are needed to finish the transmission, making the pilots to be $\mathbf{P}_i = 1$ (with $L_i = 1$). With these pilot exchanges, (5) degrades to

$$y_{j \to i} f_i - y_{i \to j} f_j = \tilde{n}. \quad (19)$$

Estimating the calibration coefficients can be performed using (15) or (16). Let us use $\mathbf{A}$ to denote $\boldsymbol{\mathcal{Y}}(\mathbf{P})^H \boldsymbol{\mathcal{Y}}(\mathbf{P})$, its element on the $i^{th}$ row and $j^{th}$ column is then given by

$$
A_{i,j} = \begin{cases} \sum_{k \ne i} |y_{k \to i}|^2 & \text{for } j = i, \\ -\, y_{j \to i}^* y_{i \to j} & \text{for } j \ne i. \end{cases}
\quad (20)
$$

Assuming a unit norm constraint, the solution given by $V_{min}(\mathbf{A})$ matches that in [22] whereas the solution under FFC corresponds to that given in [35]. Note, however, that calibration coefficients in [22], [35] are defined as the inverse of the $f_i$ in the current paper.

Other methods following the same single antenna partition scenario can be viewed as variants of the method above. For example, by allowing only the transmission between two neighboring antennas (antenna index difference is 1), (19) becomes $f_i y_{i-1 \to i} = f_{i-1} y_{i \to i-1} + \tilde{n}$. Thus, $f_i = \frac{y_{i \to i-1}}{y_{i-1 \to i}} f_{i-1} + \tilde{n}$. By setting the first antenna as the reference antenna with $f_1 = 1$, we can obtain a daisy chain calibration method as in [13], although the original was presented as a hardware-based calibration. Another variant considered in [23] consists in weighting the error metric such as $|\beta_{j \to i} f_i y_{j \to i} - \beta_{i \to j} f_j y_{i \to j}|^2$ where the weights $\beta_{j \to i}$ and $\beta_{i \to j}$ are based on
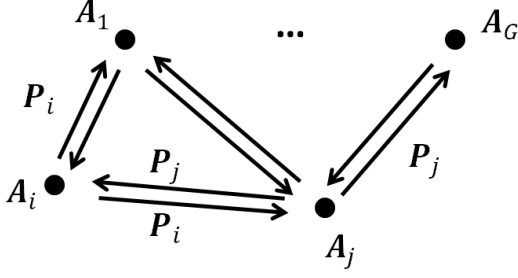
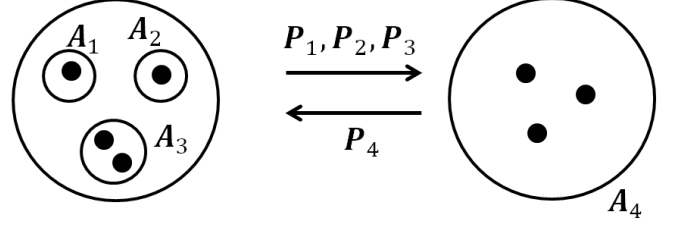Fig. 4. Method of Rogalin et al. for reciprocity calibration. Not all links between elements are plotted.



Fig. 5. Example of full Avalanche calibration with 7 antennas partitioned into 4 groups. Group 1, 2, 3 have already been calibrated, and group 4 is to be calibrated.

the SNR of the intra-array channel between antenna element $i$ and $j$.

### C. Avalanche

Avalanche [25] is a family of fast recursive calibration methods. The algorithm successively uses already calibrated parts of the antenna array to calibrate uncalibrated antennas which, once calibrated, are merged into the calibrated array. A full Avalanche calibration may be expressed under the unified framework by considering $M = \frac{1}{2}G(G-1) + 1$ antennas where $G$ is the number of groups of antennas partitioning the set of antenna elements as follows: group $A_1$ contains antenna 1, group $A_2$ contains antenna 2, group $A_3$ contains antennas 3 and 4, etc. until group $A_G$ that contains the last $G-1$ antennas. In other terms, group $A_i$ contains $M_i = \max(1, i-1)$ antennas. Moreover, in the method proposed in [25], each group $A_i$ uses $L_i = 1$ channel use by sending the pilot $\mathbf{P}_i = \mathbf{1}_{M_i \times 1}$. An example with 7 antenna elements partitioned into 4 antenna groups, where we use group 1, 2, 3 (assumed to be already calibrated) to calibrate group 4, is shown in Fig. 5. In this case, (5) becomes

$$(\mathbf{y}_{j \to i}^T * \mathbf{P}_i^T)\mathbf{f}_i - (\mathbf{P}_j^T * \mathbf{y}_{i \to j}^T)\mathbf{f}_j = \widetilde{\mathbf{n}}_{ij}. \qquad (21)$$

In [25], the authors exploited an online version of the LS estimator using previously estimated calibration parameters $\hat{\mathbf{f}}_1, \ldots, \hat{\mathbf{f}}_{i-1}$ by minimizing

$$\begin{aligned} \hat{\mathbf{f}}_i &= \arg\min_{\mathbf{f}_i} \sum_{j=1}^{i-1} \left\| (\mathbf{y}_{j \to i}^T * \mathbf{P}_i^T)\mathbf{f}_i - (\mathbf{P}_j^T * \mathbf{y}_{i \to j}^T)\hat{\mathbf{f}}_j \right\|^2 \\ &= (\mathbf{Y}_i^H \mathbf{Y}_i)^{-1} \mathbf{Y}_i^H \mathbf{a}_i, \qquad (22) \end{aligned}$$

where $\mathbf{Y}_i = \begin{bmatrix} \mathbf{y}_{1 \to i} & \mathbf{y}_{2 \to i} & \cdots & \mathbf{y}_{i-1 \to i} \end{bmatrix}^T \in \mathbb{C}^{(i-1) \times M_i}$, and $\mathbf{a}_i = [(\mathbf{P}_1^T * \mathbf{y}_{i \to 1}^T)\hat{\mathbf{f}}_1, \ldots, (\mathbf{P}_{i-1}^T * \mathbf{y}_{i \to i-1}^T)\hat{\mathbf{f}}_{i-1}] \in \mathbb{C}^{(i-1) \times 1}$. Two things should be noted, firstly, $\mathbf{f}_1, \ldots, \mathbf{f}_{i-1}$ are replaced by their estimated version which causes error propagation: estimation errors on a given calibration coefficient will propagate to subsequently calibrated antenna elements. Secondly, in order for (22) to be well-defined, i.e., in order for $\mathbf{Y}_i^H \mathbf{Y}_i$ to be invertible, it is necessary that $M_i \leqslant i - 1$. Note that this necessary condition is specific to the considered online LS estimator (22) and is more restrictive than the identifiability condition exposed in Section III-B.

## V. FAST CALIBRATION: OPTIMAL ANTENNA GROUPING

The general calibration framework in Section III opens up possibilities for new calibration schemes by using new ways to group up antennas. In this section we show that considering groups of antennas can potentially reduce the total number of channel uses necessary for calibration; we derive the theoretical limit on the smallest number of groups (and associated channel uses) needed to perform calibration.

We first address the problem of finding the smallest number of groups enabling calibration of the whole array while ensuring identifiability at each step, by finding the best choices for the $L_i$ in order to see to what extent optimizing the group based calibration can speed up the calibration process. Let us consider the case where the total number of channel uses available for calibration is fixed to $K$. We derive the number of pilot transmissions for each group, $L_1, \ldots, L_G$, that would maximize the total number of antennas that can be calibrated, i.e.,

$$\max_{(L_1, \ldots, L_G)} \left[ \sum_{j=2}^{G} \sum_{i=1}^{j-1} L_j L_i + 1 \right], \quad \text{subject to } \sum_{i=1}^{G} L_i = K. \qquad (23)$$

As shown in Appendix A, the solution of this discrete optimization problem is attained when the number of pilot transmissions for each group is equal to 1, i.e., $L_i = 1$ for any $i$ and $G = K$; note that the Avalanche approach is optimal in this sense. In this case, the number of antennas that can be calibrated is $\frac{1}{2}G(G-1) + 1$. Thus, for a given array size $M$, the number of channel uses grows only of the order of $\sqrt{M}$, which is faster than $\mathcal{O}(M)$ in Argos and the method of Rogalin et al.[5] [22]. Remark also that it is not necessary for the groups to be of equal size.

## VI. OPTIMAL ESTIMATION AND PERFORMANCE LIMITS

In order to derive estimation error bounds for the reciprocity parameters, we should not exclude a priori any data obtained during the training phase, which is what we shall assume here. In this section, we derive the CRB and associated ML

---

[5]The number of channel uses needed by the method in [22] is $M$ if we perform round-robin broadcasting for each antenna assuming that the all channels between antennas are constant during the whole calibration process whereas it would be $\mathcal{O}(M^2)$ if we perform bi-directional transmission independently for each antenna pair. Please refer to Section VII for more details.

estimation for the unified calibration scheme based on antenna partition. In order to obtain tractable results, we rely on a bilinear model to represent the calibration process. From (3), we have

$$\begin{aligned}\mathbf{Y}_{i\to j} &= \mathbf{R}_j \mathbf{C}_{i\to j}\mathbf{T}_i\mathbf{P}_i + \mathbf{N}_{i\to j}\\ &= \underbrace{\mathbf{R}_j \mathbf{C}_{i\to j}\mathbf{R}_i^T}_{\boldsymbol{\mathcal{H}}_{i\to j}} \mathbf{F}_i\mathbf{P}_i + \mathbf{N}_{i\to j},\end{aligned} \quad (24)$$

where $\mathbf{F}_i = \mathbf{R}_i^{-T}\mathbf{T}_i$ is the calibration matrix for group $i$. We define $\boldsymbol{\mathcal{H}}_{i\to j} = \mathbf{R}_j\mathbf{C}_{i\to j}\mathbf{R}_i^T$ to be a auxiliary internal channel (not corresponding to any physically measurable quantity) that appears as a nuisance parameter in the estimation of the calibration parameters. Note that the auxiliary channel $\boldsymbol{\mathcal{H}}_{i\to j}$ inherits the reciprocity from the OTA channel $\mathbf{C}_{i\to j}$: $\boldsymbol{\mathcal{H}}_{i\to j} = \boldsymbol{\mathcal{H}}_{j\to i}^T$. Upon applying the vectorization operator for each bidirectional transmission between groups $i$ and $j$, we have, similarly to (6)

$$\text{vec}(\mathbf{Y}_{i\to j}) = (\mathbf{P}_i^T * \boldsymbol{\mathcal{H}}_{i\to j})\,\mathbf{f}_i + \text{vec}(\mathbf{N}_{i\to j}). \quad (25)$$

On the reverse direction, using $\boldsymbol{\mathcal{H}}_{i\to j} = \boldsymbol{\mathcal{H}}_{j\to i}^T$, we have

$$\text{vec}(\mathbf{Y}_{j\to i}^T) = (\boldsymbol{\mathcal{H}}_{i\to j}^T * \mathbf{P}_j^T)\mathbf{f}_j + \text{vec}(\mathbf{N}_{j\to i})^T. \quad (26)$$

Alternatively, (25) and (26) may also be written as

$$\begin{cases} \text{vec}(\mathbf{Y}_{i\to j}) = [(\mathbf{F}_i\mathbf{P}_i)^T \otimes \mathbf{I}]\,\text{vec}(\boldsymbol{\mathcal{H}}_{i\to j}) + \text{vec}(\mathbf{N}_{i\to j})\\ \text{vec}(\mathbf{Y}_{j\to i}^T) = [\mathbf{I} \otimes (\mathbf{P}_j^T\mathbf{F}_j)]\,\text{vec}(\boldsymbol{\mathcal{H}}_{i\to j}) + \text{vec}(\mathbf{N}_{j\to i}). \end{cases} \quad (27)$$

Stacking these observations into a vector $\mathbf{y} = \left[\text{vec}(\mathbf{Y}_{1\to 2})^T\,\text{vec}(\mathbf{Y}_{2\to 1}^T)^T\,\text{vec}(\mathbf{Y}_{1\to 3})^T\dots\right]^T$, the above two alternative formulations can be summarized into

$$\begin{aligned}\mathbf{y} &= \boldsymbol{\mathcal{H}}(\mathbf{h},\mathbf{P})\mathbf{f} + \mathbf{n}\\ &= \boldsymbol{\mathcal{F}}(\mathbf{f},\mathbf{P})\mathbf{h} + \mathbf{n},\end{aligned} \quad (28)$$

where $\mathbf{h} = \left[\text{vec}(\boldsymbol{\mathcal{H}}_{1\to 2})^T\,\text{vec}(\boldsymbol{\mathcal{H}}_{1\to 3})^T\,\text{vec}(\boldsymbol{\mathcal{H}}_{2\to 3})^T\dots\right]^T$, and $\mathbf{n}$ is the corresponding noise vector. The composite matrices $\boldsymbol{\mathcal{H}}$ and $\boldsymbol{\mathcal{F}}$ are given by,

$$\boldsymbol{\mathcal{H}}(\mathbf{h},\mathbf{P}) = \begin{bmatrix} \mathbf{P}_1^T * \boldsymbol{\mathcal{H}}_{1\to 2} & 0 & 0 & \cdots\\ 0 & \boldsymbol{\mathcal{H}}_{1\to 2}^T * \mathbf{P}_2^T & 0 & \cdots\\ \mathbf{P}_1^T * \boldsymbol{\mathcal{H}}_{1\to 3} & 0 & 0 & \cdots\\ 0 & 0 & \boldsymbol{\mathcal{H}}_{1\to 3}^T * \mathbf{P}_3^T & \cdots\\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

$$\boldsymbol{\mathcal{F}}(\mathbf{f},\mathbf{P}) = \begin{bmatrix} \mathbf{P}_1^T\mathbf{F}_1\otimes\mathbf{I} & 0 & 0 & 0 & \cdots\\ \mathbf{I}\otimes\mathbf{P}_2^T\mathbf{F}_2 & 0 & 0 & 0 & \cdots\\ 0 & \mathbf{P}_1^T\mathbf{F}_1\otimes\mathbf{I} & 0 & 0 & \cdots\\ 0 & \mathbf{I}\otimes\mathbf{P}_3^T\mathbf{F}_3 & 0 & 0 & \cdots\\ 0 & 0 & \mathbf{P}_2^T\mathbf{F}_2\otimes\mathbf{I} & 0 & \cdots\\ 0 & 0 & \mathbf{I}\otimes\mathbf{P}_3^T\mathbf{F}_3 & 0 & \cdots\\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}. \quad (29)$$

The scenario is now identical to that encountered in some blind channel estimation scenarios and hence we can take advantage of some existing tools [36], [37], which we exploit next.

### A. Cramér-Rao bound

Treating $\mathbf{h}$ and $\mathbf{f}$ as deterministic unknown parameters, and assuming that the receiver noise $\mathbf{n}$ is distributed as $\mathcal{CN}(0,\sigma^2\mathbf{I})$, the Fisher Information Matrix (FIM) $\mathbf{J}$ for jointly estimating $\mathbf{f}$ and $\mathbf{h}$ can immediately be obtained from (28) as

$$\mathbf{J} = \frac{1}{\sigma^2}\begin{bmatrix}\boldsymbol{\mathcal{H}}^H\\ \boldsymbol{\mathcal{F}}^H\end{bmatrix}\begin{bmatrix}\boldsymbol{\mathcal{H}} & \boldsymbol{\mathcal{F}}\end{bmatrix}. \quad (30)$$

The computation of the CRB requires $\mathbf{J}$ to be non-singular. However, for the problem at hand, $\mathbf{J}$ is inherently singular. In fact, the calibration factors (and the auxiliary channel) can only be estimated up to a complex scale factor since the received data (28) involves the product of the channel and the calibration factors, $\boldsymbol{\mathcal{H}}\mathbf{f} = \boldsymbol{\mathcal{F}}\mathbf{h}$. As a result the FIM has the following null space [38], [39]

$$\mathbf{J}\begin{bmatrix}\mathbf{f}\\ -\mathbf{h}\end{bmatrix} = \frac{1}{\sigma^2}\begin{bmatrix}\boldsymbol{\mathcal{H}} & \boldsymbol{\mathcal{F}}\end{bmatrix}^H(\boldsymbol{\mathcal{H}}\mathbf{f} - \boldsymbol{\mathcal{F}}\mathbf{h}) = \mathbf{0}. \quad (31)$$

To determine the CRB when the FIM is singular, constraints have to be added to regularize the estimation problem. As the calibration parameters are complex, one complex constraint corresponds to two real constraints. Another issue is that we are mainly interested in the CRB for $\mathbf{f}$, the parameters of interest, in the presence of the nuisance parameters $\mathbf{h}$. Hence we are only interested in the $(1,1)$ block of the inverse of the $2\times 2$ block matrix $\mathbf{J}$ in (30). Incorporating the effect of the constraint (11) on $\mathbf{f}$, we can derive from [39] the following constrained CRB for $\mathbf{f}$

$$\text{CRB}_{\mathbf{f}} = \sigma^2\mathcal{V}_{\mathbf{f}}\left(\mathcal{V}_{\mathbf{f}}^H\boldsymbol{\mathcal{H}}^H\boldsymbol{\mathcal{P}}_{\boldsymbol{\mathcal{F}}}^{\perp}\boldsymbol{\mathcal{H}}\mathcal{V}_{\mathbf{f}}\right)^{-1}\mathcal{V}_{\mathbf{f}}^H \quad (32)$$

where $\boldsymbol{\mathcal{P}}_{\mathbf{X}} = \mathbf{X}(\mathbf{X}^H\mathbf{X})^\dagger\mathbf{X}^H$ and $\boldsymbol{\mathcal{P}}_{\mathbf{X}}^{\perp} = \mathbf{I} - \boldsymbol{\mathcal{P}}_{\mathbf{X}}$ are the projection operators on resp. the column space of matrix $\mathbf{X}$ and its orthogonal complement, and $\dagger$ corresponds to the Moore-Penrose pseudo inverse. Note that in some group calibration scenarios, $\boldsymbol{\mathcal{F}}^H\boldsymbol{\mathcal{F}}$ can be singular (i.e, $\mathbf{h}$ could be not identifiable even if $\mathbf{f}$ is identifiable or even known). The $M\times(M-1)$ matrix $\mathcal{V}_{\mathbf{f}}$ is such that its column space spans the orthogonal complement of that of $\frac{\partial\mathcal{C}(f)}{\partial\mathbf{f}^*}$, i.e., $\boldsymbol{\mathcal{P}}_{\mathcal{V}_{\mathbf{f}}} = \boldsymbol{\mathcal{P}}_{\frac{\partial\mathcal{C}}{\partial\mathbf{f}^*}}^{\perp}$.

It is shown in [38], [39], [40] that a choice of constraints such that their linearized version $\frac{\partial\mathcal{C}}{\partial\mathbf{f}^*}$ fills up the null space of the FIM results in the lowest CRB, while not adding information in subspaces where the data provides information. One such choice is the set (12), (13) (NPC). Another choice is (14) with $\mathbf{g} = \mathbf{f}$. With such constraints, $\frac{\partial\mathcal{C}}{\partial\mathbf{f}^*}\sim\mathbf{f}$ which spans the null space of $\boldsymbol{\mathcal{H}}^H\boldsymbol{\mathcal{P}}_{\boldsymbol{\mathcal{F}}}^{\perp}\boldsymbol{\mathcal{H}}$. The CRB then corresponds to the pseudo inverse of the FIM and (32) becomes

$$\text{CRB}_{\mathbf{f}} = \sigma^2\left(\boldsymbol{\mathcal{H}}^H\boldsymbol{\mathcal{P}}_{\boldsymbol{\mathcal{F}}}^{\perp}\boldsymbol{\mathcal{H}}\right)^\dagger. \quad (33)$$

If the FCC constraint is used instead (i.e., (14) with $\mathbf{g} = \mathbf{e}_1$, $c = 1$), the corresponding CRB is (32) where $\mathcal{V}_{\mathbf{f}}$ corresponds now to an identity matrix without the first column (and hence its column space is the orthogonal complement of that of $\mathbf{e}_1$).

Note that [24] also addresses the CRB for a scenario where transmission happens one antenna at a time. The relative calibration factors are derived from the absolute Tx and Rx side calibration parameters, which become identifiable because

a model is introduced for the internal propagation channel. In this Gaussian prior the mean is taken as the line of sight (LoS) component (distance induced delay and attenuation) and complex Gaussian non-LoS (NLOS) components are contributing to the covariance of this channel as a scaled identity matrix. The scale factor is taken 60dB below the mean channel power. This implies an almost deterministic prior for the (almost known) channel and would result in underestimation of the CRB, as noted in [24, Sec. III-E-2].

## B. Maximum likelihood estimation

We now turn our focus to the design of an optimal estimator. From (28) we get the negative log-likelihood up to an additive constant, as

$$\frac{1}{\sigma^2}\|\mathbf{y} - \mathcal{H}(\mathbf{h}, \mathbf{P})\mathbf{f}\|^2 = \frac{1}{\sigma^2}\|\mathbf{y} - \mathcal{F}(\mathbf{f}, \mathbf{P})\mathbf{h}\|^2 . \quad (34)$$

The maximum likelihood estimator of $(\mathbf{h}, \mathbf{f})$, obtained by minimizing (34), can be computed using alternating optimization on $\mathbf{h}$ and $\mathbf{f}$, which leads to a sequence of quadratic problems. As a result, for given $\mathbf{f}$, we find $\hat{\mathbf{h}} = (\mathcal{F}^H \mathcal{F})^{-1} \mathcal{F}^H \mathbf{y}$ and for given $\mathbf{h}$, we find $\hat{\mathbf{f}} = (\mathcal{H}^H \mathcal{H})^{-1} \mathcal{H}^H \mathbf{y}$. This leads to the Alternating Maximum Likelihood (AML) algorithm (Algorithm 1) [36], [37] which iteratively maximizes the likelihood by alternating between the desired parameters $\mathbf{f}$ and the nuisance parameters $\mathbf{h}$ for the formulation (28)[6].

---

**Algorithm 1** Alternating maximum likelihood (AML)

1: **Initialization:** Initialize $\hat{\mathbf{f}}$ using existing calibration methods (e.g. the method in IV-B) or as a vector of all 1's.
2: **repeat**
3:     Construct $\mathcal{F}$ as in (29) using $\hat{\mathbf{f}}$.
    $\hat{\mathbf{h}} = (\mathcal{F}^H \mathcal{F})^{-1} \mathcal{F}^H \tilde{\mathbf{y}}$
4:     Construct $\mathcal{H}$ as in (29) using $\hat{\mathbf{h}}$.
    $\hat{\mathbf{f}} = (\mathcal{H}^H \mathcal{H})^{-1} \mathcal{H}^H \tilde{\mathbf{y}}$
5: **until** the difference on the calculated $\hat{\mathbf{f}}$ between two iterations is small enough.

---

## C. Maximum likelihood vs. least squares

At first, it would seem that the ML and CRB formulations above are unrelated to the LS method introduced in Section III and used in most existing works. However, consider again the received signal in a pair $(i, j)$ as in (27). Eliminating the common auxiliary channel $\mathcal{H}_{i \to j}$, we get the elementary equation (4) for the LS method (15) or (16). Equivalently to (6), one obtains

$$\mathcal{Y}(\mathbf{P})\mathbf{f} = \mathcal{F}^{\perp H}\mathbf{y} = \tilde{\mathbf{n}}, \quad (35)$$

---

[6]The method used in [24] to derive the ML estimator, although called "Expectation Maximization" in the original paper, actually corresponds to the AML scheme, but using quadratic regularization terms for both $\mathbf{f}$ and $\mathbf{h}$ which can be interpreted as Gaussian priors and which may improve estimation in ill-conditioned cases.

where

$$\mathcal{F}^{\perp} = \begin{bmatrix} \mathbf{I} \otimes (\mathbf{F}_2 \mathbf{P}_2)^* & 0 & 0 & 0 & \dots \\ -(\mathbf{F}_1 \mathbf{P}_1)^* \otimes \mathbf{I} & 0 & 0 & 0 & \dots \\ 0 & \mathbf{I} \otimes (\mathbf{F}_3 \mathbf{P}_3)^* & 0 & 0 & \dots \\ 0 & -(\mathbf{F}_1 \mathbf{P}_1)^* \otimes \mathbf{I} & 0 & 0 & \dots \\ 0 & 0 & \mathbf{I} \otimes (\mathbf{F}_3 \mathbf{P}_3)^* & 0 & \dots \\ 0 & 0 & -(\mathbf{F}_2 \mathbf{P}_2)^* \otimes \mathbf{I} & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$
$$(36)$$

such that the column space of $\mathcal{F}^{\perp}$ corresponds to the orthogonal complement of the column space of $\mathcal{F}$ (see Appendix B) assuming that either $M_i \geq L_i$ or $L_i \geq M_i$ for all $1 \leq i \leq G$. Now, the ML criterion in (34) is separable in $\mathbf{f}$ and $\mathbf{h}$. Optimizing (34) w.r.t. $\mathbf{h}$ leads to $\mathbf{h} = (\mathcal{F}^H \mathcal{F})^{\dagger} \mathcal{F}^H \mathbf{y}$ as mentioned earlier. Substituting this estimate for $\mathbf{h}$ into (34) yields a ML estimator $\hat{\mathbf{f}}$ minimizing

$$\mathbf{y}^H \mathcal{P}_{\mathcal{F}}^{\perp} \mathbf{y} = \mathbf{y}^H \mathcal{P}_{\mathcal{F}^{\perp}} \mathbf{y} = \mathbf{y}^H \mathcal{F}^{\perp} (\mathcal{F}^{\perp H} \mathcal{F}^{\perp})^{\dagger} \mathcal{F}^{\perp H} \mathbf{y}, \quad (37)$$

where we used $\mathcal{P}_{\mathcal{F}}^{\perp} = \mathcal{P}_{\mathcal{F}^{\perp}}$. This should be compared to the least-squares method which consists in minimizing $\|\mathcal{F}^{\perp H}\mathbf{y}\|^2 = \|\mathcal{Y}\mathbf{f}\|^2$ in (15) or (16). Hence (37) can be interpreted as an optimally weighted least-squares method since from (28) $\mathcal{F}^{\perp H}\mathbf{y} = \mathcal{F}^{\perp H}\mathbf{n} = \tilde{\mathbf{n}}$ leads to colored noise with covariance matrix $\sigma^2 \mathcal{F}^{\perp H} \mathcal{F}^{\perp}$. The compressed log-likelihood in (37) can now be optimized using a variety of iterative techniques such as Iterative Quadratic ML (IQML), Denoised IQML (DIQML) or Pseudo-Quadratic ML (PQML) [37], and initialized with the least-squares method. It is not clear though whether accounting for the optimal weighting in ML would lead to significant gains in performance. The weighting matrix (before inversion) $\mathcal{F}^{\perp H} \mathcal{F}^{\perp}$ is block diagonal with a square block corresponding to the pair of antenna groups $(i, j)$ being of dimension $L_i L_j$. If all $L_i = 1$, then $\mathcal{F}^{\perp H} \mathcal{F}^{\perp}$ is a diagonal matrix. If furthermore all $M_i = 1$ (groups of isolated antennas), all pilots are of equal magnitude, and if all calibration factors would be of equal magnitude, then $\mathcal{F}^{\perp H} \mathcal{F}^{\perp}$ would be just a multipe of identity and hence would not represent any weighting. We shall leave this topic for further exploration. In any case, the fact that the CRB derived above and the ML and LS methods are all based on the signal model (28) shows that, the CRB above is the appropriate CRB for the estimation methods discussed here.

## D. Calibration bias at low SNR

Whereas the CRB applies to unbiased estimators, at low SNR the estimators are biased which turns out to lead to mean square error (MSE) saturation. In the case of a norm constraint, $\|\hat{\mathbf{f}}\|^2 = \|\mathbf{f}\|^2$, due to the triangle inequality

$$\|\hat{\mathbf{f}} - \mathbf{f}\| \leq \|\hat{\mathbf{f}}\| + \|\mathbf{f}\| = 2\|\mathbf{f}\|, \quad (38)$$

$\text{MSE} = \mathbb{E}[\|\hat{\mathbf{f}} - \mathbf{f}\|^2] \leq 4\|\mathbf{f}\|^2$. However, MSE saturation occurs also in the case of a linear constraint. We shall provide here only some brief arguments. For a linear constraint of the form (14), the least-squares method leads to (15). As the SNR decreases, the noise part $\mathcal{N}$ of $\mathcal{Y}$ will eventually dominate $\mathcal{Y}$. Hence $\hat{\mathbf{f}} = \frac{c}{\mathbf{g}^H (\mathcal{N}^H \mathcal{N})^{-1} \mathbf{g}} (\mathcal{N}^H \mathcal{N})^{-1} \mathbf{g}$ in which

the coefficients as LS estimation coefficients will tend to be bounded. To take a short-cut, consider replacing $\mathcal{N}^H\mathcal{N}$ by its mean $\mathbb{E}[\mathcal{N}^H\mathcal{N}] = c'\,\mathbf{I}$. Then we get $\hat{\mathbf{f}} = \frac{c}{\mathbf{g}^H\mathbf{g}}\mathbf{g}$ which is clearly bounded. Hence $\hat{\mathbf{f}}$ will be strongly biased with bounded MSE.

## VII. NON-COHERENT ACCUMULATION

### A. Overview

We have assumed in Sections III and IV that the channel is constant during the whole calibration process, which may become questionable if the number of antennas becomes very large since more time is then needed to accomplish the whole calibration process. As a consequence, it is possible that we cannot accumulate enough observations within a single channel coherence time and frequency block. In this section, we consider such calibration algorithms, which can jointly use data accumulated during several independent fades of the OTA channel; since the requirement to calibrate during a single coherence interval of the channel is lifted, we denote this by *non-coherent* accumulation of calibration data. Such approaches are essential for the calibration of massive MIMO systems.

Let us consider the method of Rogalin et al. as an example. If the channel is constant during the whole calibration process, one can readily use the (coherent) method detailed in Section IV-B, broadcasting pilots from each antenna in a round-robin manner when all other antennas are listening, thus $M$ slots are needed to accomplish the whole process. On the other hand, if the coherence time is not large enough, a non-coherent way to accumulate observations can be performing bi-directional transmissions for each antenna pair independently (in this case, we only require that the forward and backward transmissions are performed during the same coherence slot for each antenna pair); this requires therefore $M(M-1)$ slots. Here, we see that the non-coherent accumulation is enabled at the cost of spending more resources on calibration ($M(M-1)$ transmissions vs. $M$ transmissions for the coherent case). Some papers also implicitely use non-coherent accumulations; see for example [41] who derives a Total Least-Squares (TLS) estimator from such measurements.

Let us extend the signal model in Section III by allowing to accumulate measurements over several time slots beyond the channel coherence time (the channel can only be assumed constant within each slot, not necessarily across the slots). We assume that these are indexed by $1 \leq t \leq T$, so that $T$ represents the number of coherent slots at disposal. Clearly, the OTA reciprocity equation $\mathbf{C}_{i\to j} = \mathbf{C}_{j\to i}^T$ holds only for measurements obtained during the same time slot. However, measurements related to several groups of antennas obtained during multiple non-coherent time slots can be successfully combined to perform joint calibration of the complete array, as shown next. Let us assume that, during a given coherent slot $t$, a subset $\mathcal{G}(t)$ of the groups forming the partition of the array transmit training signals; we require that $\mathcal{G}(t)$ has at least two elements. When group $A_i$, $i \in \mathcal{G}(t)$ is transmitting, the received signal at group $A_j$, $j \in \mathcal{G}(t)$, $j \neq i$ is written

as $\mathbf{Y}_{j\to i,t} = \mathbf{R}_j\mathbf{C}_{i\to j,t}\mathbf{T}_i\mathbf{P}_{i,t} + \mathbf{N}_{j,t}$, and $\mathbf{Y}_{i\to j,t}$ is defined similarly. (5) then becomes

$$(\mathbf{Y}_{j\to i,t}^T * \mathbf{P}_{i,t}^T)\mathbf{f}_i - (\mathbf{P}_{j,t}^T * \mathbf{Y}_{i\to j,t}^T)\mathbf{f}_j = \widetilde{\mathbf{n}}_{ij,t}. \quad (39)$$

Stacking these equations similarly to (6), but with respect to the $i,j \in \mathcal{G}(t)$, gives $\boldsymbol{\mathcal{Y}}_t(\mathbf{P}_t)\mathbf{f} = \widetilde{\mathbf{n}}_t$ for each time slot $t$.

### B. LS estimation

The LS estimator of the calibration matrix is thus, taking into account all observations accumulated over the $T$ slots,

$$\hat{\mathbf{f}} = \arg\min_{\mathbf{f}} \sum_{t=1}^{T} \sum_{\substack{i,j \in \mathcal{G}(t) \\ i \neq j}} \left\| (\mathbf{Y}_{j\to i,t}^T * \mathbf{P}_{i,t}^T)\mathbf{f}_i - (\mathbf{P}_{j,t}^T * \mathbf{Y}_{i\to j,t}^T)\mathbf{f}_j \right\|^2$$
$$= \arg\min_{\mathbf{f}} \|\boldsymbol{\mathcal{Y}}(\mathbf{P})\mathbf{f}\|^2, \quad (40)$$

where the minimum is taken either under the constraint $f_1 = 1$ or $\|\mathbf{f}\| = 1$ and $\boldsymbol{\mathcal{Y}}(\mathbf{P}) = [\boldsymbol{\mathcal{Y}}_1(\mathbf{P}_1)^T, \ldots, \boldsymbol{\mathcal{Y}}_T(\mathbf{P}_T)^T]^T$. Therefore, the approach of (40) is very similar to (15) and (16). This shows that calibration using a joint estimator based on non-coherent measurements can be readily implemented by making sure that the measurements $\mathbf{Y}_{j\to i,t}$ and $\mathbf{Y}_{i\to j,t}$ appearing in each term of the sum above have been obtained during the same coherence interval. Note also that this approach can allow to collect multiple measurements across independent channel fades between the same pair $(i,j)$ of antenna groups, hence providing a way to increase the accuracy (by averaging over multiple noise realizations) and robustness (by minimizing the effect of a single catastrophic realization of the internal channel which could yield a rank-deficient set of linear equations for a given $t$) of the estimator.

### C. Optimal grouping

Statements similar to those in Section V can be made for non-coherent group-based fast calibration. The maximization proposed in Section V is still valid in this context leading to an optimal number of groups equal to the number of coherent slots $G = K$. Therefore, since $\frac{1}{2}K(K-1)$ independent rows in $\boldsymbol{\mathcal{Y}}(\mathbf{P})$ are accumulated per coherent slot, if we fix the number of antennas to $M$, the number of coherent slots $T$ should satisfy $\frac{T}{2}K(K-1) \geq M-1$ in order to calibrate all antenna elements. Note that the total number of calibrated antennas, equal to $\frac{T}{2}K(K-1) + 1$, is linear in $T$ and quadratic in $K$, which confirms that it is more valuable to perform coherent measurements in order to speed up the calibration process. However, non-coherent accumulations allow to perform measurements sparsely in time. Such a calibration process can be interleaved with the normal data transmission or reception, leading to vanishing resource overhead.

## VIII. NUMERICAL VALIDATION

In this section, we assess numerically the performance of various calibration algorithms and also compare them against their CRBs. We first evaluate the proposed group-based fast calibration method from Section V. We use MSE $= \mathbb{E}[\|\hat{\mathbf{f}} - \mathbf{f}\|^2]$ as the performance evaluation metric and the CRB as

benchmark. The Tx and Rx calibration parameters for the BS antennas are assumed to have random phases uniformly distributed over $[-\pi, \pi]$ and amplitudes uniformly distributed in the range $[1 - \delta, 1 + \delta]$ where $\delta = 0.1$. Except for the first coefficients which are fixed to 1 so that $f_1 = 1$ for the true $\mathbf{f}$. In this way, regardless of whether the FCC or the NPC (i.e. (12),(13) with $c = ||\mathbf{f}||^2$) constraints are used, direct comparison of $\hat{\mathbf{f}}$ to $\mathbf{f}$ is possible for the MSE computation (in which the expectation is replaced by sample averaging). For a fair comparison across different schemes, the number of channel uses should be the same. Hence, we compare the fast calibration method of Section V against the Avalanche scheme proposed in [25]. Note that the Argos and Rogalin methods are not fast algorithms as they need channel uses of the order of $M$, so they cannot be compared with the fast calibration methods. The number of antennas that transmit at each time instant (i.e. the group sizes of the 12 antenna groups) is shown in Table I. FC-I corresponds to a fast calibration scheme where the antenna grouping is exactly the same as that of Avalanche. However, we also try a more equally partitioned grouping of antennas in FC-II. The pilots used for transmission have unit magnitudes with uniform random phases in $[-\pi, \pi]$. The channels between all the BS antennas are assumed to be i.i.d. Rayleigh fading.

TABLE I
NUMBER OF ANTENNAS TRANSMITTING AT EACH CHANNEL USE FOR TWO FAST CALIBRATION SCHEMES.

| Scheme | Antennas transmitting per channel use. $M = 64$ | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| FC-I | 1 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 8 |
| FC-II | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 6 | 6 | 6 | 6 |
| Scheme | Antennas transmitting per channel use. $M = 67$ | | | | | | | | | | | |
| FC-I | 1 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| FC-II | 5 | 5 | 5 | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |

The performance of these schemes is depicted in Fig. 6 for $M = 64$. From Section V, it can be seen that the minimal number of channel uses required for calibration is $G = 12 = \lceil \sqrt{2M} \rceil$. The performance is averaged over 500 realizations of channel and calibration parameters. Note that the Avalanche algorithm inherently uses the FCC in its estimation process. For comparison to methods using NPC, the Avalanche estimate $\hat{\mathbf{f}}$ is then rescaled in order to satisfy the NPC constraint.

As the CRB depends on the constraint used for calibration estimation, the corresponding CRBs for these approaches are also shown. However, note that the CRB for the FC-I grouping applies to both the Avalanche method and the proposed fast calibration method (which performs least-squares (10) over all available data jointly). For each type of constraint, there are thus 3 MSE curves (Avalanche, FC-I and FC-II) and 2 CRB curves (for FC-I and FC-II). As the MSE curve is averaged over multiple channel realizations, the CRB plotted here is also an average over the CRB values corresponding to these channel realizations.

In Fig. 6, the performance of the proposed fast calibration with the FC-I grouping outperforms that of the Avalanche scheme. With $M = 64$ and $G = 12$ antenna groups, the overall system of equations is overdetermined: from (9) with $L_i = 1$, $66 = \frac{1}{2}G(G - 1) > M - 1 = 63$. This means

that the proposed fast calibration, which exploits all data jointly for the parameter estimation, has an advantage over the Avalanche method which solves exactly determined subsets of equations and hence suffers from error propagation. Also, the performance improves when the group sizes are allocated more equitably as in grouping scheme FC-II. Intuitively, the overall estimation performance of the fast calibration would be limited by the (condition number of the) largest group size and hence it is reasonable to use a grouping scheme that tries to minimize the size of the largest antenna group. These observations hold irrespective of the constraints used. Avalanche with the FCC constraint exhibits a huge MSE and hence most portions of this curve fall outside the range of Fig. 6. Note also that the MSE in some cases falls below the CRB, see for instance the MSE NPC FC-I curve at low SNRs. This is because in this SNR region the MSE saturates due to bias and the CRB is no longer applicable as explained in Section VI-D.

It is also illustrative to consider the case of $M = 67$ antennas, which is the maximum number of antennas that can be calibrated with $G = 12$ channel uses. As shown in section V, the best strategy is to divide the antennas into $G = 12$ groups and letting each group transmit exactly once ($L_i = 1$). This then results in a linear system of 66 equations (6) plus one constraint in 67 unknowns. Indeed, (9) yields $66 = \frac{1}{2}G(G-1) = M-1 = 66$. Thus, the system of equations is exactly determined by using an appropriate constraint to resolve the scale factor ambiguity. Hence, the error attained by any LS solution would be zero and the different constraints used for estimation would only lead to different scale factors in the calibration parameter estimates. So, all the solutions would be equivalent. Also, FC-I grouping leads to a block triangular structure with square diagonal blocks for the matrix $\mathcal{Y}$ defined in (7) after removing the first column. Hence, the back substitution based solution performed by Avalanche is indeed the overall LS solution with the first coefficient known constraint. Thus, in Fig. 7 where the performance of these schemes is compared for $M = 67$, we see that the curves for Avalanche and fast calibration with the FC-I grouping overlap completely. In general, this behavior would occur whenever the number of antennas corresponds to the maximum that can be calibrated with the number of channel uses (see Sec. III-B), and the antenna grouping is similar to that for FC-I. At the range of SNRs considered, the MSE is saturated and is hence far below the CRB for this grouping. In fact, only a part of the CRB for the FC-I grouping can be seen as the rest of the curve falls outside the range of the figure. Indeed, though not shown in Fig. 7, the MSE curve with the FC-I grouping only starts to overlap with the corresponding CRB curve for SNR beyond 100dB! However, it is important to note that the performance improves dramatically with a more equitable grouping of the antennas as can be seen from the curves for the FC-II grouping in the same figure.

In Fig. 8, we consider slower transmit schemes that transmit from one antenna at a time ($G = M$) and compare their MSE performance with the CRB. The MSE with FCC for Argos, Rogalin [22] and the AML method in Algorithm 1 is plotted. As expected, the Rogalin method improves over Argos by using all the bi-directional received data. AML
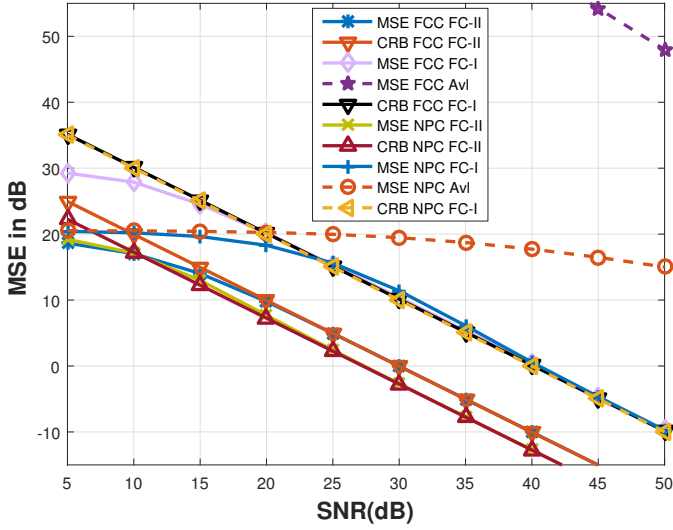
Fig. 6. Comparison of fast calibration with Avalanche scheme ($M = 64$ and the number of channel use is 12). The curves are averaged across 500 channel realizations. The performance with both the FCC and NPC constraints is shown.
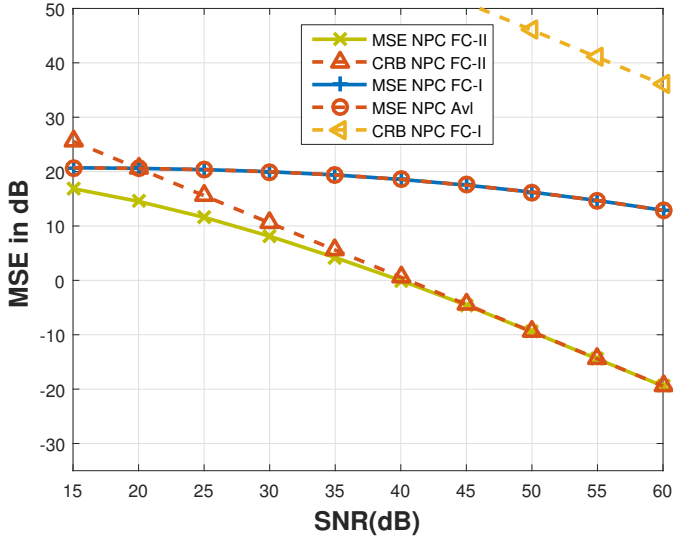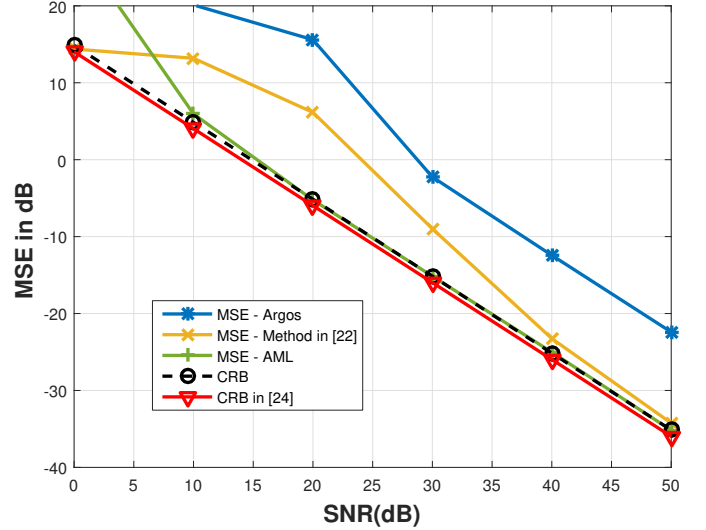


Fig. 8. Comparison of single antenna transmit schemes with the CRB ($G = M = 16$). The curves are generated over one realization of an i.i.d. Rayleigh channel and known first coefficient constraint is used.
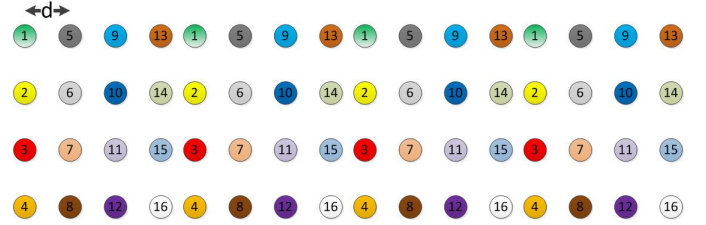


Fig. 9. 64 antennas arranged as a $4 \times 16$ grid.



Fig. 7. Comparison of fast calibration with Avalanche scheme for $M = 67$ and number of channel uses 12. The curves are averaged across 500 channel realizations. The NPC constraint is used for the MSE computation.

outperforms the Rogalin performance at low SNR. These curves are compared with the CRB derived in VI-A for the FCC case and it can be seen that the AML curve overlaps with the CRB at higher SNRs. Also plotted is the CRB as given in [24] assuming the internal propagation channel is fully known (the mean is known and the variance is negligible) and the underestimation of the MSE can be observed as expected. To bring out the difference between the two CRB derivations, the amplitude variation parameter $\delta$ is chosen to be 0.5 to increase the range of values of Tx and Rx calibration parameters.

So far, we have focused on an i.i.d. intra-array channel model and we have seen in Fig. 6 and Fig. 7 that the size of the transmission groups is an important parameter that impacts the MSE of the calibration parameter estimates. We

now consider a more realistic scenario where the intra-array channel is based on the geometry of the BS antenna array and make some observations on the choice of the antennas to form a group. We consider an array of $M = 64$ antennas arranged as in Fig. 9. The path loss $\left(4\pi \frac{d_{i \to j}}{\lambda}\right)^2$ between any two antennas $i$ and $j$ is a function of their distance $d_{i \to j}$, and $\lambda$ is the wavelength of the received signal. In the simulations, the distance between adjacent antennas, $d$, is chosen as $\frac{\lambda}{2}$. The phase of the channel between any two antennas is modeled to be a uniform random variable in $[-\pi, \pi]$. Such a model was also observed experimentally in [24]. The SNR is defined as the signal to noise ratio observed at the receive antenna nearest to the transmitter.

Continuing with the same internal channel model, consider a scenario in which antennas transmit in $G = 16$ groups of 4 each. Note that this is not the fastest grouping possible, but the example is used for the sake of illustration. We consider two different choices to form the antenna groups: 1) interleaved grouping corresponding to selecting antennas with the same numbers into one group as in Fig. 9, 2) non-interleaved grouping corresponding to selecting antennas in each column as a group. Fig. 10 shows that interleaving of the antennas results in performance gains of about 10dB. Intuitively, the interleaving of the antennas ensures that the channel from a group to the rest of the antennas is as well conditioned as possible. This example clearly shows that in addition to the
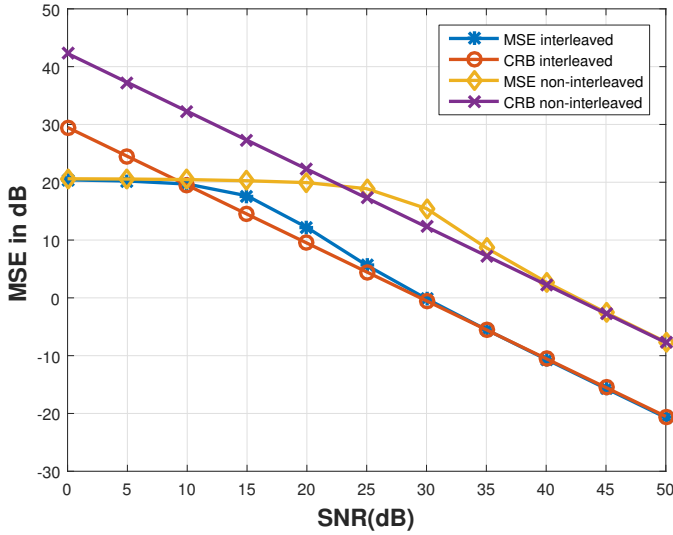
Fig. 10. Interleaved and non-interleaved MSE and CRB with NPC for an antenna transmit group size of 4 ($M = 64$ and the number of channel uses is $G = 16$).

size of the antenna groups, the choice of the antennas that go into each group also has a significant impact on the estimation quality of the calibration parameters.

## IX. CONCLUSIONS

In this work we presented an OTA calibration framework which unifies existing calibration schemes. This framework opens up new calibration possibilities. As an example, we proposed a family of fast calibration schemes based on antenna grouping. The number of channel uses needed for the whole calibration process is of the order of the square root of the antenna array size rather than scaling linearly. In fact it can be as fast as the existing Avalanche calibration method, but avoiding the severe error propagation problem, thus greatly outperforming the Avalanche method, as has been shown by simulation results.

We also presented a simple CRB formulation for the estimation of the relative calibration parameters. As the group calibration formulation encompasses the existing calibration methods, the CRB computation can be used to evaluate existing state of the art calibration methods as well. We then proposed a ML estimator and reveal the relationship between ML and LS estimation.

Moreover, we differentiated the notions of coherent and non-coherent accumulations for calibration observations. We illustrated that it is possible to perform calibration measurements using time slots that can be sparsely distributed over a relatively long time. This makes the calibration process consume a vanishing fraction of channel use resources and allows to minimize the impact on the ongoing data service.

As illustrated by simulations, for the fast calibration, interleaved grouping has a better performance than non-interleaved grouping. However, the best antenna group definitions for given antenna group sizes is still an open question. Additionally, the optimal pilot design for calibration is unknown, which is an interesting topic for future work.

## APPENDIX A
## OPTIMAL GROUPING

**Lemma 1.** *Fix* $K \geq 1$. *Let us define an optimal grouping as the solution* $G^*, L_1^*, \ldots, L_{G^*}^*$ *of*

$$\max_{\sum_{i=1}^{G} L_i = K} \sum_{i < j} L_i L_j, \tag{41}$$

*then the optimal grouping corresponds to the case* $L_1^* = \cdots = L_{G^*}^* = 1$ *with* $G^* = K$. *The number of calibrated antennas is then equal to* $\frac{1}{2} K(K - 1) + 1$.

*Proof.* Since the variables $L_1, \ldots, L_G, G$ are discrete and $\sum_{i<j} L_j L_i$ is upper bounded by $K^2$, (41) admits at least one solution. Let $\mathbf{L} = (L_1, \ldots, L_G)$ be such a solution. We reason by contradiction: suppose that there exists $j$ such that $L_j > 1$. Without loss of generality, we can suppose that $L_G > 1$. Then, we can break up group $G$ and add one group which contains a single antenna, i.e., let us consider $\mathbf{L}' = (L_1, \ldots, L_G - 1, 1)$. In that case, it holds $\sum_{i=1}^{G} L_i = \sum_{i=1}^{G+1} L_i' = K$ and

$$\sum_{j=2}^{G+1} \sum_{i=1}^{j-1} L_j' L_i'$$
$$= \sum_{j=2}^{G-1} \sum_{i=1}^{j-1} L_j' L_i' + (L_G' + L_{G+1}') \sum_{i=1}^{G-1} L_j' L_i' + L_G' L_{G+1}'$$
$$= \sum_{j=2}^{G} \sum_{i=1}^{j-1} L_j L_i + L_G' > \sum_{j=2}^{G} \sum_{i=1}^{j-1} L_j L_i$$

which contradicts the fact that $\mathbf{L}$ is solution to (41). We conclude therefore that $L_j = 1$ for any $j$ and $G^* = K$. $\square$

## APPENDIX B
## CONSTRUCTION OF $\boldsymbol{\mathcal{F}}^{\perp}$

We show in the following that the column space of $\boldsymbol{\mathcal{F}}^{\perp}$ defined by (36) spans the orthogonal complement of the column space of $\boldsymbol{\mathcal{F}}$ assuming that $\mathbf{P}_i$ is full rank for all $i$ and that either $L_i \geq M_i$ or $M_i \geq L_i$ for all i.

*Proof.* First, using $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC} \otimes \mathbf{BD})$, it holds

$$\underbrace{\left[ \mathbf{I}_{L_i} \otimes \mathbf{P}_j^T \mathbf{F}_j \quad - \mathbf{P}_i^T \mathbf{F}_i \otimes \mathbf{I}_{L_j} \right]}_{L_i L_j \times (L_i M_j + L_j M_i)} \underbrace{\begin{bmatrix} \mathbf{P}_i^T \mathbf{F}_i \otimes \mathbf{I}_{M_j} \\ \mathbf{I}_{M_i} \otimes \mathbf{P}_j^T \mathbf{F}_j \end{bmatrix}}_{(L_i M_j + L_j M_i) \times M_i M_j} = \mathbf{0} . \tag{42}$$

Then, the row space of the left matrix of (42) is orthogonal to the column space of the right matrix. As $\boldsymbol{\mathcal{F}}$ in (29) and $\boldsymbol{\mathcal{F}}^{\perp H}$ are block diagonal with blocks of the form of (42), it suffices then to prove that the following matrix $\mathbf{M}$ has full column rank, i.e., $L_i M_j + L_j M_i$, which is then also its row rank

$$\mathbf{M} := \begin{pmatrix} \mathbf{I}_{L_i} \otimes \mathbf{P}_j^T \mathbf{F}_j & -\mathbf{P}_i^T \mathbf{F}_i \otimes \mathbf{I}_{L_j} \\ (\mathbf{F}_i \mathbf{P}_i)^* \otimes \mathbf{I}_{M_j} & \mathbf{I}_{M_i} \otimes (\mathbf{F}_j \mathbf{P}_j)^* \end{pmatrix}. \tag{43}$$

Denote $\mathbf{A}_i := \mathbf{P}_i^T \mathbf{F}_i \in \mathbb{C}^{L_i \times M_i}$ and $\mathbf{A}_j := \mathbf{P}_j^T \mathbf{F}_j \in \mathbb{C}^{L_j \times M_j}$. Then, by assumption, it holds that either $\text{rank}(\mathbf{A}_i) = M_i$ and $\text{rank}(\mathbf{A}_j) = M_j$ or $\text{rank}(\mathbf{A}_i) = L_i$ and $\text{rank}(\mathbf{A}_j) =$

$L_j$. Let $\mathbf{x} = [\mathbf{x}_1^T \, \mathbf{x}_2^T]^T$ be such that $\mathbf{Mx} = 0$ and show that $\mathbf{x} = 0$. Since $\mathbf{Mx} = 0$, it holds

$$\begin{cases} (\mathbf{I}_{L_i} \otimes \mathbf{A}_j)\mathbf{x}_1 - (\mathbf{A}_i \otimes \mathbf{I}_{L_j})\mathbf{x}_2 = 0 \\ (\mathbf{A}_i^H \otimes \mathbf{I}_{M_j})\mathbf{x}_1 + (\mathbf{I}_{M_i} \otimes \mathbf{A}_j)\mathbf{x}_2 = 0. \end{cases}$$

Let $\mathbf{X}_1$ and $\mathbf{X}_2$ be matrices such that $\text{vec}(\mathbf{X}_1) = \mathbf{x}_1$ and $\text{vec}(\mathbf{X}_2) = \mathbf{x}_2$. Then

$$\begin{cases} \mathbf{A}_j \mathbf{X}_1 - \mathbf{X}_2 \mathbf{A}_i^T = 0 \\ \mathbf{X}_1 \mathbf{A}_i^* + \mathbf{A}_j^H \mathbf{X}_2 = 0. \end{cases}$$

Multiplying the first equation by $\mathbf{A}_j^H$ and the second by $\mathbf{A}_i^T$, and summing them up, we get $\mathbf{A}_j^H \mathbf{A}_j \mathbf{X}_1 + \mathbf{X}_1 (\mathbf{A}_i \mathbf{A}_i^H)^* = 0$, which is a Sylvester's equation admitting a unique solution if $\mathbf{A}_j^H \mathbf{A}_j$ and $-(\mathbf{A}_i \mathbf{A}_i^H)^*$ have no common eigenvalues. On the other hand, the eigenvalues of $\mathbf{A}_j^H \mathbf{A}_j$ and $\mathbf{A}_i \mathbf{A}_i^H$ are real positive, so common eigenvalues of $\mathbf{A}_j^H \mathbf{A}_j$ and $-(\mathbf{A}_i \mathbf{A}_i^H)^*$ can only be 0. However, this does not occur since by the assumptions either $\mathbf{A}_j^H \mathbf{A}_j$ or $\mathbf{A}_i \mathbf{A}_i^H$ is full rank. We can then conclude that $\mathbf{X}_1 = 0$, i.e., $\mathbf{x}_1 = 0$. Similarly, $\mathbf{x}_2 = 0$, which ends the proof. $\square$

## REFERENCES

[1] T. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Trans. Wireless Commun.*, vol. 9, no. 11, pp. 3590–3600, Nov. 2010.

[2] E. Larsson, O. Edfors, F. Tufvesson, and T. Marzetta, "Massive MIMO for next generation wireless systems," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 186–195, Feb. 2014.

[3] H. A. Lorentz, "het theorema van Poynting over energie in het electromagnetisch veld en een paar algemeene stellingen over de voorplanting van het licht," in *Versl. Kon. Akad. Wentensch. Amsterdam*, 1896, p. 176.

[4] G. Smith, "A direct derivation of a single-antenna reciprocity relation for the time domain," *IEEE Trans. on Antennas and Propagation*, vol. 52, no. 6, pp. 1568–1577, Jun. 2004.

[5] J. C. Guey and L. D. Larsson, "Modeling and evaluation of MIMO systems exploiting channel reciprocity in TDD mode," in *Proc. IEEE 60th Veh. Technol. Conf. (VTC)*, vol. 6, 2004, pp. 4265–4269.

[6] X. Luo, "Multi-user massive MIMO performance with calibration errors," *IEEE Trans. on Wireless Commun.*, vol. 15, no. 7, July 2016.

[7] W. Zhang, H. Ren, C. Pan, M. Chen, R. C. de Lamare, B. Du, and J. Dai, "Large-scale antenna systems with UL/DL hardware mismatch: achievable rates analysis and calibration," *IEEE Trans. on Commun.*, vol. 63, no. 4, pp. 1216–1229, 2015.

[8] X. Jiang, F. Kaltenberger, and L. Deneire, "How accurately should we calibrate a massive MIMO TDD system?" in *Proc. IEEE Intern. Conf. on Commun. (ICC) Workshops*, 2016.

[9] A. Bourdoux, B. Come, and N. Khaled, "Non-reciprocal transceivers in OFDM/SDMA systems: impact and mitigation," in *Proc. IEEE Radio and Wireless Conf. (RAWCON)*, Boston, MA, USA, Aug. 2003, pp. 183–186.

[10] K. Nishimori, K. Cho, Y. Takatori, and T. Hori, "Automatic calibration method using transmitting signals of an adaptive array for TDD systems," *Proc. IEEE Trans. on Veh. Technol.*, vol. 50, no. 6, pp. 1636–1640, 2001.

[11] K. Nishimori, T. Hiraguri, T. Ogawa, and H. Yamada, "Effectiveness of implicit beamforming using calibration technique in massive MIMO system," in *Proc. IEEE Intern. Workshop on Electromagnetics (iWEM)*, 2014, pp. 117–118.

[12] M. Petermann, M. Stefer, F. Ludwig, D. Wübben, M. Schneider, S. Paul, and K. Kammeyer, "Multi-user pre-processing in multi-antenna OFDM TDD systems with non-reciprocal transceivers," *IEEE Trans. Commun.*, vol. 61, no. 9, pp. 3781–3793, Sep. 2013.

[13] A. Benzin and G. Caire, "Internal self-calibration methods for large scale array transceiver software-defined radios," in *21th International ITG Workshop on Smart Antennas (WSA)*, Berlin, Germany, Mar. 2017.

[14] M. Guillaud, D. Slock, and R. Knopp, "A practical method for wireless channel reciprocity exploitation through relative calibration," in *Proc. Intern. Symp. Signal Process. and Its Applications*, Sydney, Australia, Aug. 2005, pp. 403–406.

[15] F. Kaltenberger, H. Jiang, M. Guillaud, and R. Knopp, "Relative channel reciprocity calibration in MIMO/TDD systems," in *Proc. Future Network and Mobile Summit*, Florence, Italy, Jun. 2010, pp. 1–10.

[16] J. Shi, Q. Luo, and M. You, "An efficient method for enhancing TDD over the air reciprocity calibration," in *Proc. IEEE Wireless Commun. and Netw. Conf.*, 2011, pp. 339–344.

[17] B. Kouassi, I. Ghauri, B. Zayen, and L. Deneire, "On the performance of calibration techniques for cognitive radio systems," in *Proc. IEEE Wireless Personal Multimedia Commun. (WPMC)*, Oct. 2011, pp. 1–5.

[18] B. Kouassi, B. Zayen, R. Knopp, F. Kaltenberger, D. Slock, I. Ghauri, F. Negro, and L. Deneire, "Design and Implementation of Spatial Interweave LTE-TDD Cognitive Radio Communication on an Experimental Platform," *IEEE Wireless Commun. Mag.*, vol. 20, no. 2, 2013.

[19] R1-091794, "Hardware calibration requirement for dual layer beamforming," Huawei, 3GPP RAN1 #57, San Francisco, USA, May 2009.

[20] R1-091752, "Performance study on Tx/Rx mismatch in LTE TDD dual-layer beamforming," Nokia, Nokia Siemens Networks, CATT, ZTE, 3GPP RAN1 #57, San Francisco, USA, May 2009.

[21] C. Shepard, N. Yu, H.and Anand, E. Li, T. Marzetta, R. Yang, and L. Zhong, "Argos: Practical many-antenna base stations," in *Proc. ACM Intern. Conf. Mobile Computing and Netw. (Mobicom)*, Istanbul, Turkey, Aug. 2012, pp. 53–64.

[22] R. Rogalin, O. Bursalioglu, H. Papadopoulos, G. Caire, A. Molisch, A. Michaloliakos, V. Balan, and K. Psounis, "Scalable synchronization and reciprocity calibration for distributed multiuser MIMO," *IEEE Trans. Wireless Commun.*, vol. 13, no. 4, pp. 1815–1831, Apr. 2014.

[23] J. Vieira, F. Rusek, and F. Tufvesson, "Reciprocity calibration methods for massive MIMO based on antenna coupling," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, Austin, USA, 2014, pp. 3708–3712.

[24] J. Vieira, F. Rusek, O. Edfors, S. Malkowsky, L. Liu, and F. Tufvesson, "Reciprocity Calibration for Massive MIMO: Proposal, Modeling and Validation," *IEEE Trans. Wireless Commun.*, May 2017.

[25] H. Papadopoulos, O. Y. Bursalioglu, and G. Caire, "Avalanche: Fast RF calibration of massive arrays," in *Proc. IEEE Global Conf. on Signal and Information Process. (GlobalSIP)*, Washington, DC, USA, Dec. 2014, pp. 607–611.

[26] J. Vieira, S. Malkowsky, Z. Nieman, K.and Miers, N. Kundargi, L. Liu, I. Wong, V. Owall, O. Edfors, and F. Tufvesson, "A flexible 100-antenna testbed for massive MIMO," in *Proc. IEEE Global Commun. Conf. (GLOBECOM) Workshops*, Austin, USA, 2014, pp. 287–293.

[27] X. Luo, "Robust large scale calibration for massive MIMO," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, San Diego, CA, USA, December 2015, pp. 1–6.

[28] H. Wei, D. Wang, H. Zhu, J. Wang, S. Sun, and X. You, "Mutual coupling calibration for multiuser massive MIMO systems," *IEEE Trans. on Wireless Commun.*, vol. 15, no. 1, pp. 606–619, 2016.

[29] G. V. Tsoulos and M. A. Beach, "Calibration and linearity issues for an adaptive antenna system," in *Proc. IEEE 47th Veh. Technol. Conf.*, vol. 3, May 1997, pp. 1597–1600.

[30] G. V. Tsoulos, J. McGeehan, and M. A. Beach, "Space division multiple access (SDMA) field trials. 2. calibration and linearity issues," *IEE Proc. Radar, Sonar and Navigation*, vol. 145, no. 1, pp. 79–84, 1998.

[31] X. Jiang, M. Čirkić, F. Kaltenberger, E. G. Larsson, L. Deneire, and R. Knopp, "MIMO-TDD reciprocity and hardware imbalances: experimental results," in *Proc. IEEE Intern. Conf. on Commun. (ICC)*, London, United Kingdom, Jun. 2015, pp. 4949–4953.

[32] C. Khatri and C. R. Rao, "Solutions to some functional equations and their applications to characterization of probability distributions," *Sankhyā: The Indian Journal of Statistics, Series A*, pp. 167–180, 1968.

[33] C. A. Balanis, *Antenna theory: analysis and design*. John Wiley & Sons, 2016.

[34] H. Wei, W. D., and X. You, "Reciprocity of mutual coupling for TDD massive MIMO systems," in *Proc. Intern. Conf. on Wireless Commun. and Signal Process. (WCSP)*, Nanjing, China, Oct. 2015, pp. 1 – 5.

[35] R. Rogalin, O. Y. Bursalioglu, H. C. Papadopoulos, G. Caire, and A. F. Molisch, "Hardware-impairment compensation for enabling distributed large-scale MIMO," in *Proc. Information Theory and Applications (ITA) Workshop, 2013*, San Diego, California, USA., Feb. 2013, pp. 1–10.

[36] E. de Carvalho and D. Slock, *Semi–Blind Methods for FIR Multichannel Estimation*. Prentice Hall, 2000, ch. 7. [Online]. Available: http://www.eurecom.fr/publication/469

[37] E. de Carvalho, S. Omar, and D. Slock, "Performance and complexity analysis of blind FIR channel identification algorithms based on deterministic maximum likelihood in SIMO systems," *Circuits, Systems, and Signal Processing*, vol. 34, no. 4, Aug. 2012. [Online]. Available: http://www.eurecom.fr/fr/publication/3790

[38] E. de Carvalho and D. Slock, "Blind and semi-blind FIR multichannel estimation: (Global) identifiability conditions," *IEEE Trans. Sig. Proc.*, Apr. 2004.

[39] ——, "Cramér-Rao bounds for blind multichannel estimation," arXiv:1710.01605 [cs.IT], 2017. [Online]. Available: https://arxiv.org/abs/1710.01605

[40] E. de Carvalho, J. Cioffi, and D. Slock, "Cramér-Rao bounds for blind multichannel estimation," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, San Francisco, CA, USA, Nov. 2000.

[41] Z. Jiang and S. Cao, "A novel TLS-based antenna reciprocity calibration scheme in TDD MIMO systems," *IEEE Commun. Letters*, vol. PP, no. 99, 2016.