# Compressive Diffusion Strategies Over Distributed Networks for Reduced Communication Load

Muhammed O. Sayin, Suleyman S. Kozat*, *Senior Member, IEEE*

*Abstract*—We study the compressive diffusion strategies over distributed networks based on the diffusion implementation and adaptive extraction of the information from the compressed diffusion data. We demonstrate that one can achieve a comparable performance with the full information exchange configurations, even if the diffused information is compressed into a scalar or a single bit. To this end, we provide a complete performance analysis for the compressive diffusion strategies. We analyze the transient, steady-state and tracking performance of the configurations in which the diffused data is compressed into a scalar or a single-bit. We propose a new adaptive combination method improving the convergence performance of the compressive diffusion strategies further. In the new method, we introduce one more freedom-of-dimension in the combination matrix and adapt it by using the conventional mixture approach in order to enhance the convergence performance for any possible combination rule used for the full diffusion configuration. We demonstrate that our theoretical analysis closely follow the ensemble averaged results in our simulations. We provide numerical examples showing the improved convergence performance with the new adaptive combination method.

*Index Terms*—Compressed diffusion, distributed network, performance analysis.

**EDICS Category: ASP-ANAL, NET-DISP**

## I. INTRODUCTION

**D**ISTRIBUTED network of nodes provides enhanced convergence performance for the applications such as source tracking, environment monitoring, and source localization [1]–[4]. In such a network, each node encounters possibly a different statistical profile, which provides broadened perspective on the monitored phenomena. In general, we would reach the best estimate with access to all observation data across the whole network since the observation of each node carries valuable information [5]. In the distributed adaptive estimation framework, we distribute the processing over the network and allow the information exchange among the nodes so that the parameter estimate of each node converges to the best estimate [4], [6].

In the distributed architectures, there are several approaches regulating the information exchange, e.g., diffusion implementation. The diffusion implementation defines a communication protocol in which only the nodes from a predefined neighborhood could exchange information with each other [1], [6]–[11]. In this framework, each node performs a local adaptive

filtering algorithm and improves its parameter estimation by fusing with the diffused parameter estimations of the neighboring nodes. The diffusion approach provides robustness against link failures and changing network topologies [6]. However, the diffusion of the parameter vector within the neighborhoods results in high amount of communication load. For example, since each node diffuses information to the neighbors, the total average number of information exchange is given by $N \times \overline{n}$ where $\overline{n}$ is the average size of a neighborhood in a network of $N$ nodes [12].

We study the compressive diffusion strategies that achieve better trade-off in terms of the amount of cooperation and the required communication load [12]. Unlike the full diffusion configuration, the compressed diffusion approach diffuses single-bit of information or a reduced dimensional data vector. The diffused data is generated through certain random projection of the local parameter estimation vector. Then, the neighboring nodes can adaptively construct the original parameter estimations based on the diffused information and fuse their individual estimates for the final estimate. This approach reduces the communication load in the spirit of the compressive sensing [12], [13]. The compression is lossy since we do not assume any sparseness or compressibility on the parameter estimation vector [13], [14]. However, the compressive diffusion approach achieves comparable convergence performance with the full diffusion configurations. Since the communication load increases far more in the large networks or highly connected network of nodes, the compressive diffusion strategies play a crucial role in achieving comparable convergence performance with significantly reduced communication load.

There exists several other approaches that reduce the communication load. In [15], within a predefined neighborhood, the parameter estimate is quantized before the diffusion in order to avoid unlimited bandwidth requirement. In [16], authors transmit the sign of the innovation sequence in the decentralized estimation framework. In [17], in a consensus network, the relative difference between the states of the nodes is exchanged by using a single bit of information. As distinct from the mentioned works, the compressive diffusion strategies substantially compress the diffused information and extract the information from the compressed data adaptively [12].

In this paper, we provide a complete performance analysis for the compressive diffusion strategies, which demonstrates comparable convergence performance of the compressed diffusion to the full information exchange configuration. We note that studying the performance of distributed networks with compressive diffusion strategies is not straight-forward

since adaptive extraction of information from the diffused data brings in an additional adaptation level. Moreover, it is rather challenging for the single-bit diffusion strategy due to the non-linear compression. However, we analyze the transient, steady-state and tracking performance of the configurations in which the diffused data is compressed into a scalar or a single-bit. We also propose a new adaptive combination method improving the performance for any conventional combination rule. In the compressive diffusion framework, we fuse the local estimates with the adaptively extracted information from substantially compressed diffusion data. The extracted information carries relatively less information than the original data. Hence, we introduce the confidence parameter concept, which adds one more freedom-of-dimension in the combination matrix. The confidence parameter determines how much we are confident with the local parameter estimation. Through the adaptation of the confidence parameter, we observe enormous enhancement in the convergence performance of the compressive diffusion strategies even for relatively long filter length.

Our main contributions include: 1) for Gaussian regressors, we analyze the transient, steady-state and tracking performance of scalar and single-bit diffusion techniques; 2) We demonstrate that our theoretical analysis accurately models the simulated results; 3) We propose a new adaptive combination method for compressive diffusion strategies, which achieves better trade-off in terms of the transient and steady state performance; 4) We provide numerical examples showing the enhanced convergence performance with the new adaptive combination method in our simulations.

We organize the paper as follows. In Section II, we explain the distributed network and diffusion implementation. In Section III, we introduce the compressive diffusion strategy, i.e., reduced-dimension and single-bit diffusion. In Section IV, we provide a global recursion model for the deviation parameters to facilitate the performance analysis. For Gaussian regressors, we analyze the mean-square convergence performance of the scalar and single-bit diffusion strategies in Section V and VI, respectively. In Section VII and VIII we analyze the steady-state and tracking performance of the scalar and single-bit diffusion approaches. In Section IX, we introduce the confidence parameter and propose a new adaptive combination method, improving the convergence performance of the compressive diffusion strategies. In Section X, we provide numerical examples demonstrating the match of theoretical and simulated results, and enhanced convergence performance with the new adaptive combination technique. We conclude the paper in Section XI with several remarks.

**Notation:** Bold lower (or upper) case letters denote the column vectors (or matrices). For a vector $\mathbf{a}$ (or matrix $\mathbf{A}$), $\mathbf{a}^T$ (or $\mathbf{A}^T$) is its ordinary transpose. $\|\cdot\|$ and $\|\cdot\|_{\mathbf{A}}$ denote the $L_2$ norm and the weighted $L_2$ norm with the matrix $\mathbf{A}$, respectively (provided that $\mathbf{A}$ is positive-definite). We work with real data for notational simplicity. For a random variable $x$ (or vector $\mathbf{x}$), $E[x]$ (or $E[\mathbf{x}]$) represents its expectation. Here, $\mathrm{Tr}(\mathbf{A})$ denotes the trace of the matrix $\mathbf{A}$. The operator $\mathrm{col}\{\cdot\}$ creates a column vector or a matrix in which the arguments of $\mathrm{col}\{\cdot\}$ locate one under the other. For a matrix argument, $\mathrm{diag}\{\cdot\}$ operator returns the diagonal of the matrix as a vector



Fig. 1: Distributed network of nodes and the neighborhood $\mathcal{N}_i$

and for a vector argument, it creates a diagonal matrix whose diagonal is the vector. The operator $\otimes$ takes the Kronecker tensor product of two matrices.

## II. DISTRIBUTED NETWORK

Consider a network of $N$ nodes where each node $i$ observes a true parameter[1] $\mathbf{w_o} \in \mathbb{R}^M$ through a linear model

$$d_{i,t} = \mathbf{w_o}^T \mathbf{u}_{i,t} + v_{i,t},$$

where $v_{i,t}$ denotes the temporally and spatially white noise. We assume that the regression vector $\mathbf{u}_{i,t} \in \mathbb{R}^M$ is spatially and temporally uncorrelated with the other regressors and the observation noise. If we know the whole temporal and spatial data overall network, we can obtain the parameter of interest $\mathbf{w_o}$ by minimizing the following global cost with respect to the parameter estimate $\mathbf{w}$:

$$J_{\mathrm{glob}}(\mathbf{w}) = \sum_{i=1}^{N} E\left[(d_{i,t} - \mathbf{w}^T \mathbf{u}_{i,t})^2\right]. \quad (1)$$

The stochastic gradient update for (1) leads to the global least-mean square (LMS) algorithm as

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mu \sum_{i=1}^{N} \mathbf{u}_{i,t}\left(d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{w}_t\right), \quad (2)$$

where $\mu > 0$ is the step size [7]. Note that (2) brings in significant communication burden by gathering the information overall network in a central processing unit. Additionally, centralized approach is not robust against the link failures and the changing network statistics [4], [6]. On the other hand, in the diffusion implementation framework, we utilize a protocol in which each node $i$ can only exchange information with nodes from its neighborhood $\mathcal{N}_i$ with the convention $i \in \mathcal{N}_i$ [6], [7]. This protocol distributes the processing to the nodes and provides tracking ability for time-varying statistical profiles [6].

Assuming the inner-node links are symmetric, we model the distributed network as an undirected graph where the nodes and the communication links correspond to its vertices and edges, respectively (See Fig. 1). In the distributed network, each node employs a local adaptation algorithm and benefits from the information diffused by the neighboring nodes in the construction of the final estimate [6]–[9]. For example, in [6], nodes *diffuse their parameter estimate* to the neighboring nodes and each node $i$ performs the LMS algorithm given as

$$\mathbf{w}_{i,t+1} = (\mathbf{I} - \mu_i \mathbf{u}_{i,t} \mathbf{u}_{i,t}^T)\boldsymbol{\phi}_{i,t} + \mu_i d_{i,t} \mathbf{u}_{i,t}, \quad (3)$$

[1]Although we assume a time invariant unknown system vector, we also provide the tracking performance analysis for certain non-stationary models later in the paper.

where $\mu_i > 0$ is the local step-size. The intermediate parameter vector $\phi_{i,t}$ is generated through

$$\phi_{i,t} = \sum_{j \in \mathcal{N}_i} \gamma_{i,j} \mathbf{w}_{j,t}$$

with $\gamma_{i,j}$'s are the combination weights such that $\sum_{j=1}^{N} \gamma_{i,j} = 1$ for all $i \in \{1, \cdots, N\}$. For a given network topology, the combination weights are determined according to certain combination rules such as uniform [18], the Metropolis [19], [20], relative-degree rules [8] or adaptive combiners [21].

We note that in (3) we could assign $\phi_{i,t}$ as the final estimate in which we adapt the local estimate through the local observation data and then we fuse with the diffused estimates to generate the final estimate. In [7], authors examine these approaches as combine-than-adapt (CTA) and adapt-than-combine (ATC) diffusion strategies, respectively. In this paper, we study the ATC diffusion strategy, however, the theoretical results hold for both the ATC and CTA cases for certain parameter changes provided later in the paper.

We emphasize that the diffusion of the parameter estimation vector also brings in high amount of communication load. In the next section, we introduce the compressive diffusion strategies enabling the adaptive construction of the required information from the reduced dimension diffusion.

### III. COMPRESSIVE DIFFUSION

We seek to estimate the parameter of interest $\mathbf{w_o}$ through the *reduced dimension information exchange* within the neighborhoods. In the compressed diffusion approach, unlike the full diffusion scheme, we diffuse a significantly reduced amount of information. The diffused information is generated by certain projection operator (a matrix $\mathcal{C}_{t+1}$ or a vector $\mathbf{c}_{t+1}$). Then, the neighboring nodes of $j$ generate an estimate $\mathbf{a}_{j,t+1}$ through the diffused information by using an adaptive estimation algorithm as explained later in the chapter [12]. We point out that the diffused information might have far smaller dimensions than the parameter estimation vector, which can reduce the communication load significantly. The constructed estimates, i.e., $\mathbf{a}_{j,t+1}$'s are linearly combined with the local parameter estimate through certain combination rules, similar to the full diffusion configuration.

Different from the full diffusion configuration, in the new framework, nodes have access to the constructed estimates $\mathbf{a}_{j,t}$. Hence, in the compressive diffusion implementation, we update according to

$$\mathbf{w}_{i,t+1} = \arg\min_{\mathbf{w}_i} \left\{ \gamma_{ii} \|\mathbf{w}_i - \mathbf{w}_{i,t}\|^2 + \sum_{j \in \mathcal{N}_i \setminus i} \gamma_{ij} \|\mathbf{w}_i - \mathbf{a}_{j,t}\|^2 \right.$$
$$\left. + \mu_i \left( d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{w}_i \right)^2 \right\} \quad (4)$$

such that in the update we also minimize the Euclidean distance between the local parameter estimation $\mathbf{w}_{i,t}$ and the constructed estimates $\mathbf{a}_{j,t}$ of the neighboring nodes. In order to simplify the optimization in (4), we can replace the loss term $(d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{w}_i)^2$ with the first order Taylor series expansion around $\mathbf{a}_{j,t}$, i.e.,

$$(d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{w}_i)^2 = \bar{e}_{i,t}(\mathbf{a}_{j,t})^2 - 2\bar{e}_{i,t}(\mathbf{a}_{j,t})\mathbf{u}_{i,t}^T(\mathbf{w}_i - \mathbf{a}_{j,t})$$
$$+ O(\|\mathbf{w}_i\|^2), \quad (5)$$



Fig. 2: CTA strategy in the compressive diffusion framework.

where we denote $\bar{e}_{i,t}(\mathbf{a}_{j,t}) \triangleq d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{a}_{j,t}$. Similarly, the first order Taylor series expansion around $\mathbf{w}_{i,t}$ leads

$$(d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{w}_i)^2 = e_{i,t}^2 - 2e_{i,t}\mathbf{u}_{i,t}^T(\mathbf{w}_i - \mathbf{w}_{i,t}) + O(\|\mathbf{w}_i\|^2), \quad (6)$$

where $e_{i,t} \triangleq d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{w}_{i,t}$. Since $\sum_{j \in \mathcal{N}_i} \gamma_{ij} = 1$, the approximations (5) and (6) in (4) yields

$$\mathbf{w}_{i,t+1} = \arg\min_{\mathbf{w}_i} \left\{ \gamma_{ii}\|\mathbf{w}_i - \mathbf{w}_{i,t}\|^2 + \sum_{j \in \mathcal{N}_i \setminus i} \gamma_{ij}\|\mathbf{w}_i - \mathbf{a}_{j,t}\|^2 \right.$$
$$+ \mu_i \gamma_{ii} \left[ e_{i,t}^2 - 2e_{i,t}\mathbf{u}_{i,t}^T(\mathbf{w}_i - \mathbf{w}_{i,t}) \right]$$
$$\left. + \mu_i \sum_{j \in \mathcal{N}_i \setminus i} \gamma_{ij} \left[ \bar{e}_{i,t}(\mathbf{a}_{j,t})^2 - 2\bar{e}_{i,t}(\mathbf{a}_{j,t})\mathbf{u}_{i,t}^T(\mathbf{w}_i - \mathbf{a}_{j,t}) \right] \right\}. \quad (7)$$

The minimized term in (7) is a convex function of $\mathbf{w}_i$ and the Hessian matrix $2\mathbf{I}_M \succ \mathbf{0}$ is positive definite. Hence, taking derivative and equating zero, we get the following update

$$\mathbf{w}_{i,t+1} = \phi_{i,t+1} + \mu_i \mathbf{u}_{i,t}(d_{i,t} - \mathbf{u}_{i,t}^T \phi_{i,t+1}), \quad (8)$$

where

$$\phi_{i,t+1} = \gamma_{ii}\mathbf{w}_{i,t} + \sum_{j \in \mathcal{N}_i \setminus i} \gamma_{ij}\mathbf{a}_{j,t}, \quad (9)$$

which is similar to the distributed LMS algorithm (3). Note that if we interchange $\phi_{i,t}$ and $\mathbf{w}_{i,t}$, in other words, when we assign the outcome of the combination as the final estimate rather than the outcome of the adaptation, we have the following algorithm:

$$\phi_{i,t+1} = \mathbf{w}_{i,t} + \mu_i \mathbf{u}_{i,t}(d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{w}_{i,t}), \quad (10)$$
$$\mathbf{w}_{i,t+1} = \gamma_{ii}\phi_{i,t+1} + \sum_{j \in \mathcal{N}_i \setminus i} \gamma_{ij}\mathbf{a}_{j,t+1}. \quad (11)$$

We point out that (8) and (9) are the CTA diffusion strategy while (10) and (11) are the ATC diffusion strategy. Fig. 2 and 3 summarize the compressive diffusion strategy for the CTA and ATC strategies where $j_k \in \mathcal{N}_i$. We next introduce different approaches to generate the diffused information (which are used to construct $a_{j,t+1}$'s).

In the compressive diffusion approach, irrespective of the final estimate we always diffuse the linear transformation of the outcome of the adaptation, e.g., we diffuse $z_{t+1} = \mathcal{C}_{t+1}^T \mathbf{w}_{i,t}$ in the CTA strategy and $z_{t+1} = \mathcal{C}_{t+1}^T \phi_{i,t+1}$ in the ATC strategy. Since we aim to use the most current parameter estimate in the construction of $\mathbf{a}_{j,t+1}$'s (since the

Fig. 3: ATC strategy in the compressive diffusion framework.

most current estimate intuitively contains more information [22]). We update according to

$$\mathbf{a}_{j,t+1} = \arg\min_{\mathbf{a}_j} \left\{ \|\mathbf{a}_j - \mathbf{a}_{j,t}\|^2 + \eta_j \|\mathbf{z}_{t+1} - \mathcal{C}_{t+1}^T \mathbf{a}_j\|^2 \right\},$$

where we choose the diffused data as the desired signal and try to minimize the mean-square of the difference between the estimate $\hat{\mathbf{z}}_{t+1} = \mathcal{C}_{t+1}^T \mathbf{a}_j$ and $\mathbf{z}_{t+1}$. The first order Taylor series approximation of the loss term $\|\mathbf{z}_{t+1} - \hat{\mathbf{z}}_{t+1}\|^2$ around $\mathbf{a}_{j,t}$ yields the following update

$$\mathbf{a}_{j,t+1} = \mathbf{a}_{j,t} + \eta_j \mathcal{C}_{t+1} (\mathbf{z}_{t+1} - \mathcal{C}_{t+1} \mathbf{a}_{j,t}) \tag{12}$$

where $\eta_j > 0$ is the construction step size. We note that in [12] the reduced dimension diffusion approach constructs $\mathbf{a}_{j,t+1}$'s through the minimum disturbance principle and resulted update involves $[\mathcal{C}_{t+1}^T \mathcal{C}_{t+1}]^{-1}$ as the normalization term. The constructed estimates $\mathbf{a}_{j,t+1}$'s are combined with the outcome of the local adaptation algorithm through (9) or (11).

We next introduce a methods where the information exchange is only a single bit [12]. When we construct $\mathbf{a}_{j,t}$ at node $j$, assuming $\mathbf{a}_{j,t}$'s are initialized with the same value, node $i \in \mathcal{N}_j$ has access to the exchanged estimate $\mathbf{a}_{j,t}$. Hence, we can perform the construction update at each neighboring node via the diffusion of the estimation error defined as

$$\epsilon_{j,t+1} \triangleq z_{t+1} - \hat{z}_{t+1}.^2$$

Note that this does not influence the communication load, however, through the access to the exchange estimate $\mathbf{a}_{j,t+1}$ we can further reduce the communication load. Using the well-known sign algorithm [5], we can construct $\mathbf{a}_{j,t+1}$ as

$$\mathbf{a}_{j,t+1} = \mathbf{a}_{j,t} + \eta_j \mathbf{c}_{t+1}\operatorname{sign}(\epsilon_{j,t+1}), \tag{13}$$

where $\mathbf{c}_{t+1}$ is the projection vector. Hence, we can repeat (13) at each neighboring node via the diffusion of $z_{j,t+1} = \operatorname{sign}(\epsilon_{j,t+1})$ only and then we combine with the local estimate by using (9) or (11).

**Remark 3.1:** The compressive diffusion strategy reduces the communication load by constructing an estimate $\mathbf{a}_{i,t+1}$ corresponding to the original estimate $\phi_{i,t+1}$ through the diffused information, i.e., the linear transformation of $\phi_{i,t+1}$ with the projection operator $\mathcal{C}_{t+1}$ or $\mathbf{c}_{t+1}$. We note that the projection operator plays crucial role in the construction algorithms (12) and (13). We choose randomized projection operator that spans the whole parameter space in order to avoid

---

²In order to facilitate the performance analyzes, we redefine $\epsilon_{j,t+1}$ in (16).

biased convergence which degrades the performance [12]. We point out that the randomized projection matrix (or vector) could be generated at each node synchronously provided that each node use the same *seed* for the pseudo-random generator mechanism [23]. Such seed exchanges and the synchronisation can be done periodically by using pilot signals without a serious increase in the communication load [24].

In the next section, we introduce a global model gathering all network operations into a single update.

## IV. GLOBAL MODEL

For a vector projection operator, we write the reduced dimension (12) and single bit (13) diffusion approaches for the ATC diffusion strategy in a compact form as

$$\phi_{i,t+1} = \mathbf{w}_{i,t} + \mu_i \mathbf{u}_{i,t} e_{i,t} \tag{14}$$

$$\mathbf{a}_{j,t+1} = \mathbf{a}_{j,t} + \eta_j \mathbf{c}_{t+1} h(\epsilon_{j,t+1}) \tag{15}$$

$$\mathbf{w}_{i,t+1} = \gamma_{i,i}\phi_{i,t+1} + \sum_{j \in \mathcal{N}_i \setminus i} \gamma_{i,j}\mathbf{a}_{j,t+1}$$

where $e_{i,t} = d_{i,t} - \mathbf{u}_{i,t}^T \mathbf{w}_{i,t}$ and $\epsilon_{j,t+1} = \mathbf{c}_{t+1}^T (\phi_{j,t+1} - \mathbf{a}_{j,t})$. For reduced dimension and single bit diffusion approaches, $h(\epsilon_{j,t+1}) = \epsilon_{j,t+1}$ and $h(\epsilon_{j,t+1}) = \operatorname{sign}(\epsilon_{j,t+1})$, respectively.

Next, we apply the following simplifications to facilitate the performance analyzes. First, we assume that at each node we use a different projection vector, e.g., for node j, we use $\mathbf{c}_{j,t}$. Second, for sufficiently small $\mu_i$, we may substitute $\phi_{i,t+1}$ with $\phi_{i,t}$ in (15) (which is justified through simulations). With that simplifications, we can rewrite the update as

$$\mathbf{a}_{j,t+1} = \mathbf{a}_{j,t} + \eta_j \mathbf{c}_{j,t} h(\epsilon_{j,t}),$$

where we redefine the construction error as

$$\epsilon_{j,t} \triangleq \mathbf{c}_{j,t}^T(\phi_{j,t} - \mathbf{a}_{j,t}). \tag{16}$$

Note that we change $\mathbf{c}_{j,t+1}$ with $\mathbf{c}_{j,t}$ to be consistent with the introduced simplification.

For the state-space representation that collects all network operations into a single update, we define the following global parameters:

$$\phi_t = \operatorname{col}\{\phi_{1,t}, \ldots, \phi_{N,t}\}, \ \mathbf{a}_t = \operatorname{col}\{\mathbf{a}_{1,t}, \ldots, \mathbf{a}_{N,t}\},$$
$$\mathbf{w}_t = \operatorname{col}\{\mathbf{w}_{1,t}, \ldots, \mathbf{w}_{N,t}\}, \ \underline{\mathbf{w}}_{\mathbf{o}} = \operatorname{col}\{\mathbf{w}_{\mathbf{o}}, \ldots, \mathbf{w}_{\mathbf{o}}\}$$

with $MN \times 1$ dimensions and

$$\mathbf{e}_t = \operatorname{col}\{e_{1,t}, \ldots, e_{N,t}\}, \ \boldsymbol{\epsilon}_t = \operatorname{col}\{\epsilon_{1,t}, \ldots, \epsilon_{N,t}\},$$
$$\mathbf{d}_t = \operatorname{col}\{d_{1,t}, \ldots, d_{N,t}\}, \ \mathbf{v}_t = \operatorname{col}\{v_{1,t}, \ldots, v_{N,t}\}$$

with $N \times 1$ dimensions. The combination matrix is given by

$$\mathbf{\Gamma} = \begin{bmatrix} \gamma_{11} & \cdots & \gamma_{1N} \\ \vdots & \ddots & \vdots \\ \gamma_{N1} & \cdots & \gamma_{NN} \end{bmatrix}$$

and we denote $\mathbf{G} \triangleq \mathbf{\Gamma} \otimes \mathbf{I}_M$. Additionally, the regression and projection vectors yields the following $MN \times N$ global matrices

$$\mathbf{U}_t \triangleq \begin{bmatrix} \mathbf{u}_{1,t} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \mathbf{u}_{N,t} \end{bmatrix}, \ \mathbf{C}_t \triangleq \begin{bmatrix} \mathbf{c}_{1,t} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \mathbf{c}_{N,t} \end{bmatrix}.$$

Indeed, we can model the network with compressive diffusion strategy as a larger network in which each node $i$ has an imaginary counterpart which diffuses $\mathbf{a}_{i,t}$ to the neighbors of $i$, which is similar to the full diffusion configuration. The real nodes only get information from the imaginary nodes and do not diffuse any information. In that case, the network can be modelled as a directed graph with asymmetric inner node links and the combination matrix is given by

$$\tilde{\boldsymbol{\Gamma}} = \begin{bmatrix} \boldsymbol{\Gamma}_D & \boldsymbol{\Gamma}_C \\ \mathbf{0} & \mathbf{I} \end{bmatrix},$$

where $\boldsymbol{\Gamma}_D = \mathrm{diag}\{\boldsymbol{\Gamma}\}$ and $\boldsymbol{\Gamma}_C = \boldsymbol{\Gamma} - \boldsymbol{\Gamma}_D$. Then, we can write $\mathbf{w}_t$ in terms of $\boldsymbol{\phi}_t$ and $\mathbf{a}_t$ as

$$\mathbf{w}_t = \mathbf{G_D}\boldsymbol{\phi}_t + \mathbf{G_C}\mathbf{a}_t, \tag{17}$$

where $\mathbf{G_D} \triangleq \boldsymbol{\Gamma}_D \otimes \mathbf{I}_M$ and $\mathbf{G_C} \triangleq \boldsymbol{\Gamma}_C \otimes \mathbf{I}_M$. The state-space representation is given by

$$\boldsymbol{\phi}_{t+1} = \mathbf{w}_t + \mathbf{M}\mathbf{U}_t\mathbf{e}_t,$$
$$\mathbf{a}_{t+1} = \mathbf{a}_t + \mathbf{N}\mathbf{C}_t\mathbf{h}(\boldsymbol{\epsilon}_t),$$
$$\mathbf{w}_{t+1} = \mathbf{G_D}\boldsymbol{\phi}_{t+1} + \mathbf{G_C}\mathbf{a}_{t+1},$$

where

$$\mathbf{M} \triangleq \mathrm{diag}\{[\mu_1,\ldots,\mu_N]\} \otimes \mathbf{I}_M,$$

$$\mathbf{N} \triangleq \mathrm{diag}\{[\eta_1,\ldots,\eta_N]\} \otimes \mathbf{I}_M$$

and $\mathbf{h}(\boldsymbol{\epsilon}_t) = \mathrm{col}\{h(\epsilon_{1,t}),\cdots,h(\epsilon_{N,t})\}$. We obtain the global deviation vectors as

$$\tilde{\boldsymbol{\phi}}_t \triangleq \underline{\mathbf{w}_o} - \boldsymbol{\phi}_t \quad \text{and} \quad \tilde{\mathbf{a}}_t \triangleq \underline{\mathbf{w}_o} - \mathbf{a}_t. \tag{18}$$

Since $\boldsymbol{\Gamma}\underline{\mathbf{1}} = \underline{\mathbf{1}}$,

$$\mathbf{G}\underline{\mathbf{w}_o} = \underline{\mathbf{w}_o} \tag{19}$$

then the global deviation update yields

$$\tilde{\boldsymbol{\phi}}_{t+1} = \mathbf{G_D}\tilde{\boldsymbol{\phi}}_t + \mathbf{G_C}\tilde{\mathbf{a}}_t - \mathbf{M}\mathbf{U}_t\mathbf{e}_t, \tag{20}$$
$$\tilde{\mathbf{a}}_{t+1} = \tilde{\mathbf{a}}_t - \mathbf{N}\mathbf{C}_t\mathbf{h}(\boldsymbol{\epsilon}_t). \tag{21}$$

In (22), we represent the global deviation updates (20) and (21) in a single equation or equivalently

$$\tilde{\boldsymbol{\psi}}_{t+1} = \mathbf{X}\tilde{\boldsymbol{\psi}}_t - \mathbf{D}\mathbf{Y}_t\underline{\mathbf{h}}(\mathbf{e}_t,\boldsymbol{\epsilon}_t), \tag{23}$$

where $\tilde{\boldsymbol{\psi}}_t \triangleq \mathrm{col}\{\tilde{\boldsymbol{\phi}}_t, \tilde{\mathbf{a}}_t\}$. Based on the weighted-energy recursion of (23), in the next sections, we analyze the mean-square convergence performance of scalar and single-bit diffusion approaches separately for Gaussian regressors.

## V. Scalar Diffusion with Gaussian Regressors

For the one-dimension diffusion approach, (23) yields

$$\tilde{\boldsymbol{\psi}}_{t+1} = \mathbf{X}\tilde{\boldsymbol{\psi}}_t - \mathbf{D}\mathbf{Y}_t\underline{\mathbf{e}}_t, \tag{24}$$

where $\underline{\mathbf{e}}_t \triangleq \mathrm{col}\{\mathbf{e}_t,\boldsymbol{\epsilon}_t\}$. By (17), (18) and (19), we note that $\mathbf{e}_t$ is given by

$$\mathbf{e}_t = \mathbf{U}_t^T(\mathbf{G_D}\tilde{\boldsymbol{\phi}}_t + \mathbf{G_C}\tilde{\mathbf{a}}_t) + \mathbf{v}_t. \tag{25}$$

Similarly, we have

$$\boldsymbol{\epsilon}_t = \mathbf{C}_t^T(-\tilde{\boldsymbol{\phi}}_t + \tilde{\mathbf{a}}_t). \tag{26}$$

Hence, through (25) and (26), we obtain the global estimation error $\underline{\mathbf{e}}_t$ as

$$\underline{\mathbf{e}}_t = \begin{bmatrix} \mathbf{U}_t & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_t \end{bmatrix}^T \underbrace{\begin{bmatrix} \mathbf{G_D} & \mathbf{G_C} \\ -\mathbf{I} & \mathbf{I} \end{bmatrix}}_{\mathbf{Z}} \begin{bmatrix} \tilde{\boldsymbol{\phi}}_t \\ \tilde{\mathbf{a}}_t \end{bmatrix} + \underbrace{\begin{bmatrix} \mathbf{v}_t \\ \mathbf{0} \end{bmatrix}}_{\mathbf{n}_t}$$
$$= \mathbf{Y}_t^T\mathbf{Z}\tilde{\boldsymbol{\psi}}_t + \mathbf{n}_t. \tag{27}$$

Through (27), we rewrite (24) as

$$\tilde{\boldsymbol{\psi}}_{t+1} = \mathbf{X}\tilde{\boldsymbol{\psi}}_t - \mathbf{D}\mathbf{Y}_t(\mathbf{Y}_t^T\mathbf{Z}\tilde{\boldsymbol{\psi}}_t + \mathbf{n}_t)$$
$$= (\mathbf{X} - \mathbf{D}\mathbf{Y}_t\mathbf{Y}_t^T\mathbf{Z})\tilde{\boldsymbol{\psi}}_t - \mathbf{D}\mathbf{Y}_t\mathbf{n}_t. \tag{28}$$

We utilize the weighted-energy relation relating the energy of the error and deviation quantities in the performance analyzes through a weighting matrix $\boldsymbol{\Sigma}$. Then, we obtain

$$\tilde{\boldsymbol{\psi}}_{t+1}^T\boldsymbol{\Sigma}\tilde{\boldsymbol{\psi}}_{t+1} = [(\mathbf{X} - \mathbf{D}\mathbf{Y}_t\mathbf{Y}_t^T\mathbf{Z})\tilde{\boldsymbol{\psi}}_t - \mathbf{D}\mathbf{Y}_t\mathbf{n}_t]^T\boldsymbol{\Sigma}$$
$$\times [(\mathbf{X} - \mathbf{D}\mathbf{Y}_t\mathbf{Y}_t^T\mathbf{Z})\tilde{\boldsymbol{\psi}}_t - \mathbf{D}\mathbf{Y}_t\mathbf{n}_t]$$
$$= \tilde{\boldsymbol{\psi}}_t^T(\mathbf{X} - \mathbf{D}\mathbf{Y}_t\mathbf{Y}_t^T\mathbf{Z})^T\boldsymbol{\Sigma}(\mathbf{X} - \mathbf{D}\mathbf{Y}_t\mathbf{Y}_t^T\mathbf{Z})\tilde{\boldsymbol{\psi}}_t$$
$$- 2\mathbf{n}_t^T\mathbf{Y}_t^T\mathbf{D}\boldsymbol{\Sigma}(\mathbf{X} - \mathbf{D}\mathbf{Y}_t\mathbf{Y}_t^T\mathbf{Z})\tilde{\boldsymbol{\psi}}_t$$
$$+ \mathbf{n}_t^T\mathbf{Y}_t^T\mathbf{D}\boldsymbol{\Sigma}\mathbf{D}\mathbf{Y}_t\mathbf{n}_t.$$

Since we assume the observation noise $\mathbf{v}_t$ is independent from the network statistics, the weighted energy relation for (28) is given by

$$E\|\tilde{\boldsymbol{\psi}}_{t+1}\|_{\boldsymbol{\Sigma}}^2 = E\|\tilde{\boldsymbol{\psi}}_t\|_{\boldsymbol{\Sigma}'}^2 + E[\mathbf{n}_t^T\mathbf{Y}_t^T\mathbf{D}\boldsymbol{\Sigma}\mathbf{D}\mathbf{Y}_t\mathbf{n}_t] \tag{29}$$

where

$$\boldsymbol{\Sigma}' \triangleq (\mathbf{X} - \mathbf{D}\mathbf{Y}_t\mathbf{Y}_t^T\mathbf{Z})^T\boldsymbol{\Sigma}(\mathbf{X} - \mathbf{D}\mathbf{Y}_t\mathbf{Y}_t^T\mathbf{Z})$$
$$= \mathbf{X}^T\boldsymbol{\Sigma}\mathbf{X} - \mathbf{Z}^T\mathbf{Y}_t\mathbf{Y}_t^T\mathbf{D}\boldsymbol{\Sigma}\mathbf{X} - \mathbf{X}^T\boldsymbol{\Sigma}\mathbf{D}\mathbf{Y}_t\mathbf{Y}_t^T\mathbf{Z}$$
$$+ \mathbf{Z}^T\mathbf{Y}_t\mathbf{Y}_t^T\mathbf{D}\boldsymbol{\Sigma}\mathbf{D}\mathbf{Y}_t\mathbf{Y}_t^T\mathbf{Z}.$$

Apart form the weighting matrix $\boldsymbol{\Sigma}$, $\boldsymbol{\Sigma}'$ is a random due to the data dependence. We assume the spatial and temporal independence of the regression data $\mathbf{u}_{i,t}$ and $\mathbf{c}_{j,t}$ so that $\mathbf{Y}_t$ is independent of $\tilde{\boldsymbol{\psi}}_t$. Through that assumption we can replace $\boldsymbol{\Sigma}'$ by its mean value, i.e., $\boldsymbol{\Sigma}' = E[\boldsymbol{\Sigma}']$ [5], [6]. Hence, the weighting matrix is given by

$$\boldsymbol{\Sigma}' = \mathbf{X}^T\boldsymbol{\Sigma}\mathbf{X} - \mathbf{Z}^T E\left[\mathbf{Y}_t\mathbf{Y}_t^T\right]\mathbf{D}\boldsymbol{\Sigma}\mathbf{X} - \mathbf{X}^T\boldsymbol{\Sigma}\mathbf{D}E\left[\mathbf{Y}_t\mathbf{Y}_t^T\right]\mathbf{Z}$$
$$+ \mathbf{Z}^T\mathbf{D}E\left[\mathbf{Y}_t\mathbf{Y}_t^T\boldsymbol{\Sigma}\mathbf{Y}_t\mathbf{Y}_t^T\right]\mathbf{D}\mathbf{Z}. \tag{30}$$

Note that in the last term of right hand side (RHS) of (30) we take $\mathbf{D}$'s out of the expectation thanks to the block diagonal structure of $\mathbf{D}$ and $\mathbf{Y}_t\mathbf{Y}_t^T$.

In order to calculate certain data moments in (29) and (30), we assume spatially and temporally i.i.d. Gaussian regression data such that

$$\boldsymbol{\Lambda_u} \triangleq E[\mathbf{U}_t\mathbf{U}_t^T] = \mathrm{diag}\{[\sigma_{u,1}^2, \sigma_{u,2}^2, \ldots, \sigma_{u,N}^2]\} \otimes \mathbf{I}_M$$
$$\boldsymbol{\Lambda_c} \triangleq E[\mathbf{C}_t\mathbf{C}_t^T] = \mathrm{diag}\{[\sigma_{c,1}^2, \sigma_{c,2}^2, \ldots, \sigma_{c,N}^2]\} \otimes \mathbf{I}_M.$$

Then, we obtain

$$\boldsymbol{\Lambda} \triangleq E[\mathbf{Y}_t\mathbf{Y}_t^T] = \begin{bmatrix} \boldsymbol{\Lambda_u} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Lambda_c} \end{bmatrix}.$$

In the performance analysis, convenient vectorisation notation is used to exploit the diagonal structure of matrices [5], [25]. In (29), (30), matrices have block diagonal structures,

$$\overbrace{\begin{bmatrix} \tilde{\phi}_{t+1} \\ \tilde{\mathbf{a}}_{t+1} \end{bmatrix}}^{\tilde{\psi}_{t+1}} = \overbrace{\begin{bmatrix} \mathbf{G_D} & \mathbf{G_C} \\ \mathbf{0} & \mathbf{I}_{MN} \end{bmatrix}}^{\mathbf{X}} \overbrace{\begin{bmatrix} \tilde{\phi}_t \\ \tilde{\mathbf{a}}_t \end{bmatrix}}^{\tilde{\psi}_t} - \overbrace{\begin{bmatrix} \mathbf{M} & \mathbf{0} \\ \mathbf{0} & \mathbf{N} \end{bmatrix}}^{\mathbf{D}} \overbrace{\begin{bmatrix} \mathbf{U}_t & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_t \end{bmatrix}}^{\mathbf{Y}_t} \overbrace{\begin{bmatrix} \mathbf{e}_t \\ \mathbf{h}(\epsilon_t) \end{bmatrix}}^{\underline{\mathbf{h}}(\mathbf{e}_t, \epsilon_t)} \tag{22}$$

thus, we use the block vectorisation operator bvec$\{\cdot\}$ [6]. Given an $NM \times NM$ block matrix

$$\mathbf{\Sigma} = \begin{bmatrix} \mathbf{\Sigma}_{11} & \ldots & \mathbf{\Sigma}_{1N} \\ \vdots & \ddots & \vdots \\ \mathbf{\Sigma}_{N1} & \ldots & \mathbf{\Sigma}_{NN} \end{bmatrix}$$

where each block $\mathbf{\Sigma}_{ij}$ is a $M \times M$ block. Let $\boldsymbol{\sigma}_{ij} = \text{vec}\{\mathbf{\Sigma}_{ij}\}$ with standard vec$\{\cdot\}$ operator and $\boldsymbol{\sigma}_j = \text{col}\{\boldsymbol{\sigma}_{1j}, \boldsymbol{\sigma}_{2j}, \ldots, \boldsymbol{\sigma}_{Nj}\}$, then

$$\text{bvec}\{\mathbf{\Sigma}\} = \boldsymbol{\sigma} = \text{col}\{\boldsymbol{\sigma}_1, \boldsymbol{\sigma}_2, \ldots, \boldsymbol{\sigma}_N\}. \tag{31}$$

We also use the *block Kronecker product* of two block matrices $\mathbf{A}$ and $\mathbf{B}$, denoted by $\mathbf{A} \odot \mathbf{B}$. The $ij$-block is given by

$$[\mathbf{A} \odot \mathbf{B}]_{ij} = \begin{bmatrix} \mathbf{A}_{ij} \otimes \mathbf{B}_{11} & \ldots & \mathbf{A}_{ij} \otimes \mathbf{B}_{1N} \\ \vdots & \ddots \vdots \\ \mathbf{A}_{ij} \otimes \mathbf{B}_{N1} & \ldots & \mathbf{A}_{ij} \otimes \mathbf{B}_{NN} \end{bmatrix}. \tag{32}$$

The block vectorisation operator bvec$\{\cdot\}$ (31) and the block Kronecker product (32) are related by

$$\text{bvec}\{\mathbf{A}\mathbf{\Sigma}\mathbf{B}\} = (\mathbf{B}^T \odot \mathbf{A})\boldsymbol{\sigma} \tag{33}$$

and

$$\text{Tr}\{\mathbf{A}^T \mathbf{B}\} = (\text{bvec}\{\mathbf{A}\})^T \text{bvec}\{\mathbf{B}\}. \tag{34}$$

The term in the RHS of (29) yields

$$E\left[\mathbf{n}_t^T \mathbf{Y}_t^T \mathbf{D}\mathbf{\Sigma}\mathbf{D}\mathbf{Y}_t \mathbf{n}_t\right] = \text{Tr}\left(\mathbf{\Lambda}\mathbf{D}^2 E\left[\mathbf{n}_t \mathbf{n}_t^T\right] \mathbf{\Sigma}\right)$$

and let

$$E\left[\mathbf{n}_t \mathbf{n}_t^T\right] = \mathbf{R_n} = \begin{bmatrix} \mathbf{R_v} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix},$$

where $\mathbf{R_v} \triangleq \text{diag}\{\sigma_{v,1}^2, \ldots, \sigma_{v,N}^2\} \otimes \mathbf{I}_M$. Then by (34),

$$E\left[\mathbf{n}_t^T \mathbf{Y}_t^T \mathbf{D}\mathbf{\Sigma}\mathbf{D}\mathbf{Y}_t \mathbf{n}_t\right] = \mathbf{b}^T \boldsymbol{\sigma},$$

where

$$\mathbf{b} \triangleq \text{bvec}\{\mathbf{R_n}\mathbf{D}^2\mathbf{\Lambda}\}. \tag{35}$$

The fourth-order moment in (30) yields

$$\mathbf{A} = E\left[\mathbf{Y}_t \mathbf{Y}_t^T \mathbf{\Sigma}\mathbf{Y}_t \mathbf{Y}_t^T\right],$$

where the $M \times M$ block is given by

$$[\mathbf{A}]_{ij} = \begin{cases} 2\mathbf{\Lambda}_i \mathbf{\Sigma}_{ii} \mathbf{\Lambda}_i + \mathbf{\Lambda}_i \text{Tr}\left(\mathbf{\Sigma}_{ii}\mathbf{\Lambda}_i\right) & i = j \\ \mathbf{\Lambda}_i \mathbf{\Sigma}_{ij} \mathbf{\Lambda}_j & i \neq j \end{cases}$$

thanks to the spatial and temporal independence of the regression data [5]. We note that $\mathbf{\Lambda}$ could be denoted as $\mathbf{\Lambda} = \text{diag}\{\mathbf{\Lambda}_1, \cdots, \mathbf{\Lambda}_N\}$ where $\mathbf{\Lambda}_i$ for $i = \{1, 2, \ldots, N\}$ is $M \times M$ block matrix, e.g., $\mathbf{\Lambda}_1 = \sigma_{u,1}^2 \mathbf{I}_M$. The $M \times M$ $ij$th block of $\mathbf{\Sigma}$ is denoted by $\mathbf{\Sigma}_{ij}$. Through (32), (34), we obtain

$$\text{bvec}\{\mathbf{A}\} = \mathcal{A}\boldsymbol{\sigma}$$

with $\mathcal{A} = \text{diag}\{\mathcal{A}_1, \ldots, \mathcal{A}_N\}$, $\mathcal{A}_j = \text{diag}\{\mathcal{A}_{1j}, \ldots, \mathcal{A}_{Nj}\}$ and

$$\mathcal{A}_{ij} = \begin{cases} 2\mathbf{\Lambda}_i \otimes \mathbf{\Lambda}_i + \lambda_i \lambda_i^T & i = j \\ \mathbf{\Lambda}_i \otimes \mathbf{\Lambda}_j & i \neq j \end{cases}$$

where $\lambda_i = \text{vec}\{\mathbf{\Lambda}_i\}$.

Hence, the block vectorization of the weighting matrix $\mathbf{\Sigma}'$ (30) yields

$$\text{bvec}\{\mathbf{\Sigma}'\} = \Big(\mathbf{X}^T \odot \mathbf{X}^T - (\mathbf{X}^T \odot \mathbf{Z}^T)(\mathbf{I}_{2MN} \odot \mathbf{\Lambda}\mathbf{D})$$
$$- (\mathbf{Z}^T \odot \mathbf{X}^T)(\mathbf{\Lambda}\mathbf{D} \odot \mathbf{I}_{2MN})$$
$$+ (\mathbf{Z}^T \odot \mathbf{Z}^T)(\mathbf{D} \odot \mathbf{D})\mathcal{A}\Big)\boldsymbol{\sigma}.$$

For notational simplicity, we change the weighted-norm notation such that $\|\tilde{\phi}_t\|_{\boldsymbol{\sigma}}^2$ refers to $\|\tilde{\phi}_t\|_{\mathbf{\Sigma}}^2$ where $\boldsymbol{\sigma} = \text{bvec}\{\mathbf{\Sigma}\}$. As a result, we obtain the weighted-energy recursion as

$$E\|\tilde{\psi}_{t+1}\|_{\boldsymbol{\sigma}}^2 = E\|\tilde{\psi}_t\|_{\mathbf{F}\boldsymbol{\sigma}}^2 + \mathbf{b}^T \boldsymbol{\sigma} \tag{36}$$

$$\mathbf{F} \triangleq \mathbf{X}^T \odot \mathbf{X}^T + (\mathbf{Z}^T \odot \mathbf{Z}^T)(\mathbf{D} \odot \mathbf{D})\mathcal{A}$$
$$- (\mathbf{X}^T \odot \mathbf{Z}^T)(\mathbf{I}_{2MN} \odot \mathbf{\Lambda}\mathbf{D})$$
$$- (\mathbf{Z}^T \odot \mathbf{X}^T)(\mathbf{\Lambda}\mathbf{D} \odot \mathbf{I}_{2MN}). \tag{37}$$

Through (36) and (37), we can analyze the learning, convergence and stability behavior of the network. Iterating the weighted-energy recursion, we obtain

$$E\|\tilde{\psi}_{t+1}\|_{\boldsymbol{\sigma}}^2 = E\|\tilde{\psi}_t\|_{\mathbf{F}\boldsymbol{\sigma}}^2 + \mathbf{b}^T \boldsymbol{\sigma}$$
$$E\|\tilde{\psi}_t\|_{\mathbf{F}\boldsymbol{\sigma}}^2 = E\|\tilde{\psi}_{t-1}\|_{\mathbf{F}^2\boldsymbol{\sigma}}^2 + \mathbf{b}^T \mathbf{F}\boldsymbol{\sigma}$$
$$\vdots$$
$$E\|\tilde{\psi}_1\|_{\mathbf{F}^t\boldsymbol{\sigma}}^2 = E\|\tilde{\psi}_0\|_{\mathbf{F}^{t+1}\boldsymbol{\sigma}}^2 + \mathbf{b}^T \mathbf{F}^t \boldsymbol{\sigma}.$$

Assuming the parameter estimates $\phi_{i,t}$ and $\mathbf{a}_{i,t}$ are initialized with zeros, $E\|\tilde{\psi}_0\|^2 = \|\underline{\mathbf{w}}_{\mathbf{o}}\|^2$ where $\underline{\mathbf{w}}_{\mathbf{o}} \triangleq \text{col}\{\underline{\mathbf{w}}_{\mathbf{o}}, \underline{\mathbf{w}}_{\mathbf{o}}\}$. The iterations yield

$$E\|\tilde{\psi}_{t+1}\|_{\boldsymbol{\sigma}}^2 = \|\underline{\mathbf{w}}_{\mathbf{o}}\|_{\mathbf{F}^{t+1}\boldsymbol{\sigma}}^2 + \mathbf{b}^T \left(\sum_{k=0}^t \mathbf{F}^k\right)\boldsymbol{\sigma}. \tag{38}$$

By (38), we reach the following final recursion:

$$E\|\tilde{\psi}_{t+1}\|_{\boldsymbol{\sigma}}^2 = E\|\tilde{\psi}_t\|_{\boldsymbol{\sigma}}^2 + \mathbf{b}^T \mathbf{F}^t \boldsymbol{\sigma} - \|\underline{\mathbf{w}}_{\mathbf{o}}\|_{\mathbf{F}^t(\mathbf{I}-\mathbf{F})\boldsymbol{\sigma}}^2. \tag{39}$$

**Remark 5.1:** We note that (39) is of essence since through the weighting matrix $\mathbf{\Sigma}$ we can extract information about the learning and convergence behavior of the network. In Table I, we tabulate the initial conditions (we assume the initial parameter vectors are set to $\mathbf{0}$) and the weighting matrices corresponding to various conventional performance measures.

**Remark 5.2:** In this paper, (39) provides a recursion for the weighted deviation parameter where we assign $\phi_{i,t}$ as the

final estimate instead of $\mathbf{w}_{i,t}$, which implies the CTA strategy, however, the recursion also provides the performance of the ATC strategy with appropriate combination matrix $\boldsymbol{\Sigma}$ and the initial condition (See Table I).

Next, we analyze the mean-square convergence performance of the single-bit diffusion approach for Gaussian regressors.

## VI. SINGLE-BIT DIFFUSION WITH GAUSSIAN REGRESSORS

The weighted-energy relation of (23) yields

$$
\begin{aligned}
E\left[\tilde{\boldsymbol{\psi}}_{t+1}^T \boldsymbol{\Sigma} \tilde{\boldsymbol{\psi}}_{t+1}\right] = & \; E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X}^T \boldsymbol{\Sigma} \mathbf{X} \tilde{\boldsymbol{\psi}}_t\right] \\
& - E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X}^T \boldsymbol{\Sigma} \mathbf{D} \mathbf{Y}_t \underline{\mathbf{h}}(\mathbf{e}_t, \boldsymbol{\epsilon}_t)\right] \\
& - E\left[\underline{\mathbf{h}}^T(\mathbf{e}_t, \boldsymbol{\epsilon}_t) \mathbf{Y}_t^T \mathbf{D} \boldsymbol{\Sigma} \mathbf{X} \tilde{\boldsymbol{\psi}}_t\right] \\
& + E\left[\underline{\mathbf{h}}^T(\mathbf{e}_t, \boldsymbol{\epsilon}_t) \mathbf{Y}_t^T \mathbf{D} \boldsymbol{\Sigma} \mathbf{D} \mathbf{Y}_t \underline{\mathbf{h}}(\mathbf{e}_t, \boldsymbol{\epsilon}_t)\right].
\end{aligned}
$$
(40)

We evaluate RHS of (40) term by term in order to find the variance relation. Firstly, we partition the weighting matrix as follows:

$$
\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_1 & \boldsymbol{\Sigma}_2 \\ \boldsymbol{\Sigma}_3 & \boldsymbol{\Sigma}_4 \end{bmatrix}.
$$
(41)

Through the partitioning (41), we obtain

$$
\begin{aligned}
E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X}^T \boldsymbol{\Sigma} \mathbf{D} \mathbf{Y}_t \underline{\mathbf{h}}(\mathbf{e}_t, \boldsymbol{\epsilon}_t)\right] = & \; E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X_u}^T \boldsymbol{\Sigma}_1 \mathbf{M} \mathbf{U}_t \mathbf{U}_t^T \mathbf{Z_u} \tilde{\boldsymbol{\psi}}_t\right] \\
& + E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X_u}^T \boldsymbol{\Sigma}_2 \mathbf{N} \mathbf{C}_t \mathrm{sign}\left(\mathbf{C}_t^T \mathbf{Z_d} \tilde{\boldsymbol{\psi}}_t\right)\right] \\
& + E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X_d}^T \boldsymbol{\Sigma}_3 \mathbf{M} \mathbf{U}_t \mathbf{U}_t^T \mathbf{Z_u} \tilde{\boldsymbol{\psi}}_t\right] \\
& + E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X_d}^T \boldsymbol{\Sigma}_4 \mathbf{N} \mathbf{C}_t \mathrm{sign}\left(\mathbf{C}_t^T \mathbf{Z_d} \tilde{\boldsymbol{\psi}}_t\right)\right],
\end{aligned}
$$
(42)

where we partition $\mathbf{X}$ such that $\mathbf{X} = \mathrm{col}\{\mathbf{X_u}, \mathbf{X_d}\}$. We note that the second and fourth terms in the RHS of (42) include the nonlinear $\mathrm{sign}(\cdot)$ function. It is not straight forward to evaluate the expectations with nonlinearity, thus, we introduce the following lemma.

**Lemma 1:** *Under the assumption that step-sizes are sufficiently small and the filter is sufficiently long [5], the Price's theorem leads to*

$$
\begin{aligned}
& E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X_u}^T \boldsymbol{\Sigma}_2 \mathbf{N} \mathbf{C}_t \mathrm{sign}\left(\mathbf{C}_t^T \mathbf{Z_d} \tilde{\boldsymbol{\psi}}_t\right)\right] \\
& \qquad = E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X_u}^T \boldsymbol{\Sigma}_2 \mathbf{N} \boldsymbol{\Omega}_t \mathbf{C}_t \mathbf{C}_t^T \mathbf{Z_d} \tilde{\boldsymbol{\psi}}_t\right],
\end{aligned}
$$
(43)

$$
\begin{aligned}
& E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X_d}^T \boldsymbol{\Sigma}_4 \mathbf{N} \mathbf{C}_t \mathrm{sign}\left(\mathbf{C}_t^T \mathbf{Z_d} \tilde{\boldsymbol{\psi}}_t\right)\right] \\
& \qquad = E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X_d}^T \boldsymbol{\Sigma}_4 \mathbf{N} \boldsymbol{\Omega}_t \mathbf{C}_t \mathbf{C}_t^T \mathbf{Z_d} \tilde{\boldsymbol{\psi}}_t\right],
\end{aligned}
$$
(44)

*where $\boldsymbol{\Omega}_t$ is defined as*

$$
\boldsymbol{\Omega}_t \triangleq \begin{bmatrix} \frac{E|\epsilon_{1,t}|}{E[\epsilon_{1,t}^2]} \mathbf{I}_M & \cdots & \mathbf{0}_M \\ \vdots & \ddots & \vdots \\ \mathbf{0}_M & \cdots & \frac{E|\epsilon_{N,t}|}{E[\epsilon_{N,t}^2]} \mathbf{I}_M \end{bmatrix}
$$

*Proof:* The proof is given in Appendix A. $\square$

By (42), (43), (44), the second term on the RHS of (40) is given by

$$
\begin{aligned}
& E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X}^T \boldsymbol{\Sigma} \mathbf{D} \mathbf{Y}_t \underline{\mathbf{h}}(\mathbf{e}_t, \boldsymbol{\epsilon}_t)\right] \\
& \qquad = E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X}^T \boldsymbol{\Sigma} \mathbf{D} \underline{\boldsymbol{\Omega}}_t \mathbf{Y}_t \mathbf{Y}_t^T \mathbf{Z} \tilde{\boldsymbol{\psi}}_t\right],
\end{aligned}
$$
(45)

where $\underline{\boldsymbol{\Omega}}_t$ denotes

$$
\underline{\boldsymbol{\Omega}}_t \triangleq \begin{bmatrix} \mathbf{I}_{MN} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Omega}_t \end{bmatrix}.
$$

Similarly, the third term on the RHS of (40) is evaluated as

$$
\begin{aligned}
& E\left[\underline{\mathbf{h}}^T(\mathbf{e}_t, \boldsymbol{\epsilon}_t) \mathbf{Y}_t^T \mathbf{D} \boldsymbol{\Sigma} \mathbf{X} \tilde{\boldsymbol{\psi}}_t\right] \\
& \qquad = E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{Z}^T \mathbf{Y}_t \mathbf{Y}_t^T \underline{\boldsymbol{\Omega}}_t \mathbf{D} \boldsymbol{\Sigma} \mathbf{X} \tilde{\boldsymbol{\psi}}_t\right].
\end{aligned}
$$
(46)

Through partitioning, the last term on the RHS of (40) leads to

$$
\begin{aligned}
& E\left[\underline{\mathbf{h}}^T(\mathbf{e}_t, \boldsymbol{\epsilon}_t) \mathbf{Y}_t^T \mathbf{D} \boldsymbol{\Sigma} \mathbf{D} \mathbf{Y}_t \underline{\mathbf{h}}(\mathbf{e}_t, \boldsymbol{\epsilon}_t)\right] \\
& \qquad = E\left[\mathbf{e}_t^T \mathbf{U}_t^T \mathbf{M} \boldsymbol{\Sigma}_1 \mathbf{M} \mathbf{U}_t \mathbf{e}_t\right] \\
& \qquad + E\left[\mathbf{e}_t^T \mathbf{U}_t^T \mathbf{M} \boldsymbol{\Sigma}_2 \mathbf{N} \mathbf{C}_t \mathrm{sign}(\boldsymbol{\epsilon}_t)\right] \\
& \qquad + E\left[\mathrm{sign}(\boldsymbol{\epsilon}_t)^T \mathbf{C}_t^T \mathbf{N} \boldsymbol{\Sigma}_3 \mathbf{M} \mathbf{U}_t \mathbf{e}_t\right] \\
& \qquad + E\left[\mathrm{sign}(\boldsymbol{\epsilon}_t)^T \mathbf{C}_t^T \mathbf{N} \boldsymbol{\Sigma}_4 \mathbf{N} \mathbf{C}_t \mathrm{sign}(\boldsymbol{\epsilon}_t)\right].
\end{aligned}
$$

**Corollary 1:** *Since $\mathbf{U}_t$ and $\mathbf{C}_t$ are independent from each other, similar to the Lemma 1, we obtain*

$$
\begin{aligned}
& E\left[\underline{\mathbf{h}}^T(\mathbf{e}_t, \boldsymbol{\epsilon}_t) \mathbf{Y}_t^T \mathbf{D} \boldsymbol{\Sigma} \mathbf{D} \mathbf{Y}_t \underline{\mathbf{h}}(\mathbf{e}_t, \boldsymbol{\epsilon}_t)\right] \\
& \qquad = E\left[\mathbf{e}_t^T \mathbf{U}_t^T \mathbf{M} \boldsymbol{\Sigma}_1 \mathbf{M} \mathbf{U}_t \mathbf{e}_t\right] \\
& \qquad + E\left[\mathbf{e}_t^T \mathbf{U}_t^T \mathbf{M} \boldsymbol{\Sigma}_2 \mathbf{N} \boldsymbol{\Omega}_t \mathbf{C}_t \boldsymbol{\epsilon}_t\right] \\
& \qquad + E\left[\boldsymbol{\epsilon}_t^T \mathbf{C}_t^T \boldsymbol{\Omega}_t \mathbf{N} \boldsymbol{\Sigma}_3 \mathbf{M} \mathbf{U}_t \mathbf{e}_t\right] \\
& \qquad + E\left[\mathrm{sign}(\boldsymbol{\epsilon}_t)^T \mathbf{C}_t^T \mathbf{N} \boldsymbol{\Sigma}_4 \mathbf{N} \mathbf{C}_t \mathrm{sign}(\boldsymbol{\epsilon}_t)\right].
\end{aligned}
$$
(47)

Because of the independence of the observation noise from the regression data, the first term on the RHS of (47) yields

$$
\begin{aligned}
E\left[\mathbf{e}_t^T \mathbf{U}_t^T \mathbf{M} \boldsymbol{\Sigma}_1 \mathbf{M} \mathbf{U}_t \mathbf{e}_t\right] = & \; E\left[\mathbf{v}_t^T \mathbf{U}_t^T \mathbf{M} \boldsymbol{\Sigma}_1 \mathbf{M} \mathbf{U}_t \mathbf{v}_t\right] \\
& + E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{Z_u}^T \mathbf{U}_t \mathbf{U}_t^T \mathbf{M} \boldsymbol{\Sigma}_1 \mathbf{M} \mathbf{U}_t \mathbf{U}_t^T \mathbf{Z}_u \tilde{\boldsymbol{\psi}}_t\right].
\end{aligned}
$$
(48)

For the last term on the RHS of (47), we introduce the following lemma.

**Lemma 2:** *Through the Price's theorem, we obtain*

$$
\begin{aligned}
& E\left[\mathrm{sign}(\boldsymbol{\epsilon}_t)^T \mathbf{C}_t^T \mathbf{N} \boldsymbol{\Sigma}_4 \mathbf{N} \mathbf{C}_t \mathrm{sign}(\boldsymbol{\epsilon}_t)\right] \\
& \qquad = E\left[\tilde{\boldsymbol{\psi}}_t^T \mathbf{Z_d}^T \mathbf{C}_t \mathbf{C}_t^T \mathbf{N} \boldsymbol{\Omega}_t \boldsymbol{\Sigma}_4^C \boldsymbol{\Omega}_t \mathbf{N} \mathbf{C}_t \mathbf{C}_t^T \mathbf{Z_d} \tilde{\boldsymbol{\psi}}_t\right] \\
& \qquad + E\left[\mathbf{1}^T \mathbf{C}_t^T \mathbf{N} \boldsymbol{\Sigma}_4^D \mathbf{N} \mathbf{C}_t \mathbf{1}\right],
\end{aligned}
$$
(49)

TABLE I: Initial conditions and weighting matrices for different configurations.

| Framework | $E\|\tilde{\psi}_t\|^2_{\mathbf{\Sigma}}$ | $E\|\tilde{\psi}_0\|^2_{\mathbf{\Sigma}}$ | $\mathbf{\Sigma}$ | $E\|\tilde{\psi}_t\|^2_{\mathbf{\Sigma}}$ | $E\|\tilde{\psi}_0\|^2_{\mathbf{\Sigma}}$ | $\mathbf{\Sigma}$ |
|---|---|---|---|---|---|---|
| CTA | $\frac{1}{N}E\|\tilde{\phi}_t\|^2$ | $\frac{1}{N}\|\mathbf{w_o}\|^2$ | $\frac{1}{N}\begin{bmatrix}\mathbf{I}_{MN} & \mathbf{0} \\ \mathbf{0} & \mathbf{0}\end{bmatrix}$ | $\frac{1}{N}E\|\tilde{\phi}_t\|^2_{\mathbf{\Lambda_u}}$ | $\frac{1}{N}\|\mathbf{w_o}\|^2_{\mathbf{\Lambda_u}}$ | $\frac{1}{N}\begin{bmatrix}\mathbf{\Lambda}_u & \mathbf{0} \\ \mathbf{0} & \mathbf{0}\end{bmatrix}$ |
| ATC | $\frac{1}{N}E\|\tilde{\mathbf{w}}_t\|^2$ | $\frac{1}{N}\|\mathbf{w_o}\|^2$ | $\frac{1}{N}\begin{bmatrix}\mathbf{G_D}^T\mathbf{G_D} & \mathbf{G_D}^T\mathbf{G_C} \\ \mathbf{G_C}^T\mathbf{G_D} & \mathbf{G_C}^T\mathbf{G_C}\end{bmatrix}$ | $\frac{1}{N}E\|\tilde{\mathbf{w}}_t\|^2_{\mathbf{\Lambda_u}}$ | $\frac{1}{N}\|\mathbf{w_o}\|^2_{\mathbf{\Lambda_u}}$ | $\frac{1}{N}\begin{bmatrix}\mathbf{G_D}^T\mathbf{\Lambda}_u\mathbf{G_D} & \mathbf{G_D}^T\mathbf{\Lambda}_u\mathbf{G_C} \\ \mathbf{G_C}^T\mathbf{\Lambda}_u\mathbf{G_D} & \mathbf{G_C}^T\mathbf{\Lambda}_u\mathbf{G_C}\end{bmatrix}$ |

*where $\mathbf{\Sigma}_4^D$ is the block diagonal matrix of $\mathbf{\Sigma}_4$ such that*

$$\mathbf{\Sigma}_4^D = \begin{bmatrix}\mathbf{\Theta}_{11} & \cdots & \mathbf{0}_M \\ \vdots & \ddots & \vdots \\ \mathbf{0}_M & \cdots & \mathbf{\Theta}_{NN}\end{bmatrix}$$

*with $\mathbf{\Theta}_{ii}$ is the $ii$'th $M \times M$ block of $\mathbf{\Sigma}_4$ and $\mathbf{\Sigma}_4^C = \mathbf{\Sigma}_4 - \mathbf{\Sigma}_4^D$.*

*Proof:* The proof is given in Appendix B. $\square$

As a result, by (45), (46), (47), (48) and (49); (40) leads to

$$E\|\tilde{\psi}_{t+1}\|^2_{\mathbf{\Sigma}} = E\|\tilde{\psi}_t\|^2_{\mathbf{\Sigma}'} + E\left[\mathbf{v}_t^T\mathbf{U}_t^T\mathbf{M}\mathbf{\Sigma}_1\mathbf{M}\mathbf{U}_t\mathbf{v}_t\right]$$
$$+ E\left[\mathbf{1}^T\mathbf{C}_t^T\mathbf{N}\mathbf{\Sigma}_4^D\mathbf{N}\mathbf{C}_t\mathbf{1}\right] \quad (50)$$

and

$$\mathbf{\Sigma}' = \mathbf{X}^T\mathbf{\Sigma}\mathbf{X} - \mathbf{X}^T\mathbf{\Sigma}\mathbf{D}\underline{\mathbf{\Omega}}_t\mathbf{Y}_t\mathbf{Y}_t^T\mathbf{Z} - \mathbf{Z}^T\mathbf{Y}_t\mathbf{Y}_t^T\underline{\mathbf{\Omega}}_t\mathbf{D}\mathbf{\Sigma}\mathbf{X}$$
$$+ \mathbf{Z}^T\mathbf{D}\underline{\mathbf{\Omega}}_t\mathbf{Y}_t\mathbf{Y}_t^T\tilde{\mathbf{\Sigma}}\mathbf{Y}_t\mathbf{Y}_t^T\underline{\mathbf{\Omega}}_t\mathbf{D}\mathbf{Z},$$

where $\tilde{\mathbf{\Sigma}}$ denotes

$$\tilde{\mathbf{\Sigma}} = \begin{bmatrix}\mathbf{\Sigma}_1 & \mathbf{\Sigma}_2 \\ \mathbf{\Sigma}_3 & \mathbf{\Sigma}_4^C\end{bmatrix}.$$

We again note that under the assumption that the regression data is spatially and temporally independent, we get $\mathbf{\Sigma}' = E[\mathbf{\Sigma}']$ which results

$$\mathbf{\Sigma}' = \mathbf{X}^T\mathbf{\Sigma}\mathbf{X} - \mathbf{X}^T\mathbf{\Sigma}\mathbf{D}\underline{\mathbf{\Omega}}_t\mathbf{\Lambda}\mathbf{Z} - \mathbf{Z}^T\mathbf{\Lambda}\underline{\mathbf{\Omega}}_t\mathbf{D}\mathbf{\Sigma}\mathbf{X}$$
$$+ \mathbf{Z}^T\mathbf{D}\underline{\mathbf{\Omega}}_t E\left[\mathbf{Y}_t\mathbf{Y}_t^T\tilde{\mathbf{\Sigma}}\mathbf{Y}_t\mathbf{Y}_t^T\right]\underline{\mathbf{\Omega}}_t\mathbf{D}\mathbf{Z} \quad (51)$$

and denote $\mathbf{B} \triangleq E\left[\mathbf{Y}_t\mathbf{Y}_t^T\tilde{\mathbf{\Sigma}}\mathbf{Y}_t\mathbf{Y}_t^T\right]$. Now, we resort to the vector notation, i.e., the block vectorisation operator $\text{bvec}\{\cdot\}$ and the block Kronecker product. Hence, the block vectorization of the weighting matrix $\mathbf{\Sigma}'$ (51) yields

$$\text{bvec}\{\mathbf{\Sigma}'\} = \left(\mathbf{X}^T \odot \mathbf{X}^T - (\mathbf{X}^T \odot \mathbf{Z}^T)(\mathbf{I}_{2MN} \odot \mathbf{\Lambda}\mathbf{D}\underline{\mathbf{\Omega}}_t)\right.$$
$$\left. - (\mathbf{Z}^T \odot \mathbf{X}^T)(\mathbf{\Lambda}\mathbf{D}\underline{\mathbf{\Omega}}_t \odot \mathbf{I}_{2MN})\right)\boldsymbol{\sigma}$$
$$+ (\mathbf{Z}^T \odot \mathbf{Z}^T)(\mathbf{D} \odot \mathbf{D})(\underline{\mathbf{\Omega}}_t \odot \underline{\mathbf{\Omega}}_t)\text{bvec}\{\mathbf{B}\}. \quad (52)$$

Block vectorisation of the matrix $\mathbf{B}$ is given by

$$\text{bvec}\{\mathbf{B}\} = \mathcal{A}\,\text{bvec}\{\tilde{\mathbf{\Sigma}}\}.$$

In order to denote $\text{bvec}\{\tilde{\mathbf{\Sigma}}\}$ in terms of $\boldsymbol{\sigma}$, we introduce the following matrices:

$$\mathbf{K}_1 \triangleq \text{col}\{\mathbf{0}_{MN}, \mathbf{I}_{MN}\},$$
$$\mathbf{K}_2 \triangleq \text{col}\{\mathbf{I}_{MN}, \mathbf{0}_{MN}\},$$
$$\mathbf{T_k} \triangleq \text{diag}\{\mathbf{0}_{(k-1)M}, \mathbf{I}_M, \mathbf{0}_{(N-k)M}\}.$$

We get $\mathbf{\Sigma}_4^D$ and $\tilde{\mathbf{\Sigma}}$ as

$$\mathbf{\Sigma}_4^D = \sum_{k=1}^{N}\mathbf{T_k}\mathbf{K}_2^T\mathbf{\Sigma}\mathbf{K}_2\mathbf{T_k}, \quad (53)$$
$$\tilde{\mathbf{\Sigma}} = \mathbf{\Sigma} - \mathbf{K}_2\mathbf{\Sigma}_4^D\mathbf{K}_2^T. \quad (54)$$

By (53) and (54), we obtain

$$\text{bvec}\{\tilde{\mathbf{\Sigma}}\} = \underbrace{\left(\mathbf{I} - (\mathbf{K}_2 \odot \mathbf{K}_2)\sum_{k=1}^{N}(\mathbf{T_k} \odot \mathbf{T_k})(\mathbf{K_2}^T \odot \mathbf{K}_2^T)\right)}_{\mathbf{K}}\boldsymbol{\sigma}$$
$$= \mathbf{K}\boldsymbol{\sigma}. \quad (55)$$

The $\tilde{\psi}$-free terms in (50) are evaluated as

$$E\left[\mathbf{v}_t^T\mathbf{U}_t^T\mathbf{M}\mathbf{\Sigma}_1\mathbf{M}\mathbf{U}_t\mathbf{v}_t\right] = \mathbf{b}_1^T(\mathbf{K}_1^T \odot \mathbf{K}_1^T)\boldsymbol{\sigma}, \quad (56)$$
$$E\left[\mathbf{1}^T\mathbf{C}_t^T\mathbf{N}\mathbf{\Sigma}_4^D\mathbf{N}\mathbf{C}_t\mathbf{1}\right] = \mathbf{b}_2^T(\mathbf{K}_2^T \odot \mathbf{K}_2^T)\boldsymbol{\sigma}, \quad (57)$$

where $\mathbf{b}_1 \triangleq \text{bvec}\{\mathbf{R_v}\mathbf{M}^2\mathbf{\Lambda}_u\}$ and $\mathbf{b}_2 \triangleq \text{bvec}\{\mathbf{1}\mathbf{1}^T\mathbf{N}^2\mathbf{\Lambda_c}\}$.

As a result, by (52), (55), (56) and (57), the weighted-energy relation is given by

$$E\|\tilde{\psi}_{t+1}\|^2_{\boldsymbol{\sigma}} = E\|\tilde{\psi}_t\|^2_{\mathbf{F}_t\boldsymbol{\sigma}} + \mathbf{b}^T\boldsymbol{\sigma} \quad (58)$$
$$\mathbf{F}_t = \mathbf{X}^T \odot \mathbf{X}^T - (\mathbf{X}^T \odot \mathbf{Z}^T)(\mathbf{I}_{2MN} \odot \mathbf{\Lambda}\mathbf{D}\underline{\mathbf{\Omega}}_t)$$
$$- (\mathbf{Z}^T \odot \mathbf{X}^T)(\mathbf{\Lambda}\mathbf{D}\underline{\mathbf{\Omega}}_t \odot \mathbf{I}_{2MN})$$
$$+ (\mathbf{Z}^T \odot \mathbf{Z}^T)(\mathbf{D} \odot \mathbf{D})(\underline{\mathbf{\Omega}}_t \odot \underline{\mathbf{\Omega}}_t)\mathcal{A}\mathbf{K} \quad (59)$$
$$\mathbf{b} = (\mathbf{K}_1^T \odot \mathbf{K}_1^T)^T\mathbf{b}_1 + (\mathbf{K}_2^T \odot \mathbf{K}_2^T)^T\mathbf{b}_2. \quad (60)$$

Iterating the weighted-energy recursion (58), (59) and (60), we obtain

$$E\|\tilde{\psi}_{t+1}\|^2_{\boldsymbol{\sigma}} = E\|\tilde{\psi}_t\|^2_{\mathbf{F}_t\boldsymbol{\sigma}} + \mathbf{b}^T\boldsymbol{\sigma}$$
$$E\|\tilde{\psi}_t\|^2_{\mathbf{F}_t\boldsymbol{\sigma}} = E\|\tilde{\psi}_{t-1}\|^2_{\mathbf{F}_{t-1}\mathbf{F}_t\boldsymbol{\sigma}} + \mathbf{b}^T\mathbf{F}_t\boldsymbol{\sigma}$$
$$\vdots$$
$$E\|\tilde{\psi}_1\|^2_{\mathbf{F}_1\ldots\mathbf{F}_t\boldsymbol{\sigma}} = E\|\tilde{\psi}_0\|^2_{\mathbf{F}_0\ldots\mathbf{F}_t\boldsymbol{\sigma}} + \mathbf{b}^T\mathbf{F}_1\ldots\mathbf{F}_t\boldsymbol{\sigma}.$$

In this part of the analyzes, we do not assume that the parameter vectors are initialized with zeros since such an assumption results in infinite terms in the $\mathbf{\Omega}_t$ matrix. Hence, we initialize $\mathbf{a}_t$ with $\zeta\mathbf{1}_{MN\times 1}$ where $\zeta$ has a small value (See Table II).

The iterations yield

$$E\|\tilde{\psi}_{t+1}\|^2_{\boldsymbol{\sigma}} = \|\tilde{\psi}_0\|^2_{\mathbf{\Pi}_t\boldsymbol{\sigma}} + \mathbf{b}^T\mathbf{\Delta}_t\boldsymbol{\sigma}, \quad (61)$$
$$E\|\tilde{\psi}_t\|^2_{\boldsymbol{\sigma}} = \|\tilde{\psi}_0\|^2_{\mathbf{\Pi}_{t-1}\boldsymbol{\sigma}} + \mathbf{b}^T\mathbf{\Delta}_{t-1}\boldsymbol{\sigma}, \quad (62)$$

TABLE II: Initial conditions and weighting matrices for the performance measure of the construction update for the single-bit diffusion approach (for the scalar diffusion approach, set $\zeta = 0$) and the global MSD of the ATC diffusion strategy for the single-bit diffusion approach (for the scalar diffusion approach, see Table I).

| $E\|\tilde{\psi}_t\|^2_{\underline{\Sigma}}$ | $E\|\tilde{\psi}_0\|^2_{\underline{\Sigma}}$ | $\Sigma$ |
|---|---|---|
| $\frac{1}{N}E\|\tilde{\mathbf{a}}_t\|^2$ | $\frac{1}{N}\|\underline{\mathbf{w}}_{\mathbf{o}} - \zeta\mathbf{1}\|^2$ | $\begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \frac{1}{N}\mathbf{I}_{MN} \end{bmatrix}$ |
| $\sigma^2_{\boldsymbol{\epsilon}_t} = E[\boldsymbol{\epsilon}^T_t\boldsymbol{\epsilon}_t]$ | $\zeta\mathbf{1}^T\boldsymbol{\Lambda}_{\mathbf{c}}\mathbf{1}$ | $\begin{bmatrix} \boldsymbol{\Lambda}_{\mathbf{c}} & -\boldsymbol{\Lambda}_{\mathbf{c}} \\ -\boldsymbol{\Lambda}_{\mathbf{c}} & \boldsymbol{\Lambda}_{\mathbf{c}} \end{bmatrix}$ |
| $\frac{1}{N}E\|\tilde{\mathbf{w}}_t\|^2$ | $\frac{1}{N}\|\underline{\mathbf{w}}_{\mathbf{o}} - \zeta\mathbf{G}_{\mathbf{C}}\mathbf{1}\|^2$ | $\frac{1}{N}\begin{bmatrix} \mathbf{G}_{\mathbf{D}}{}^T\mathbf{G}_{\mathbf{D}} & \mathbf{G}_{\mathbf{D}}{}^T\mathbf{G}_{\mathbf{C}} \\ \mathbf{G}_{\mathbf{C}}{}^T\mathbf{G}_{\mathbf{D}} & \mathbf{G}_{\mathbf{C}}{}^T\mathbf{G}_{\mathbf{C}} \end{bmatrix}$ |

where $\boldsymbol{\Pi}_t \overset{\triangle}{=} \prod_{i=0}^{t}\mathbf{F}_i$ and $\boldsymbol{\Delta}_t \overset{\triangle}{=} \mathbf{I} + \mathbf{F}_t + \mathbf{F}_{t-1}\mathbf{F}_t + \cdots + \mathbf{F}_1\dots\mathbf{F}_t$. We note that $\boldsymbol{\Pi}_t = \boldsymbol{\Pi}_{t-1}\mathbf{F}_t$ and $\boldsymbol{\Delta}_t = \boldsymbol{\Delta}_{t-1}\mathbf{F}_t + \mathbf{I}$. By (61) and (62), we have the following recursion

$$E\|\tilde{\psi}_{t+1}\|^2_{\boldsymbol{\sigma}} = E\|\tilde{\psi}_t\|^2_{\boldsymbol{\sigma}} - \|\tilde{\psi}_0\|^2_{\boldsymbol{\Pi}_{t-1}(\mathbf{I}-\mathbf{F}_t)\boldsymbol{\sigma}} + \mathbf{b}^T\left(\mathbf{I} - \boldsymbol{\Delta}_{t-1}(\mathbf{I} - \mathbf{F}_t)\right)\boldsymbol{\sigma}. \quad (63)$$

We point out that $\boldsymbol{\Pi}_{-1} = \mathbf{I}_{(2MN)^2}$ and $\boldsymbol{\Delta}_{-1} = \mathbf{0}_{(2MN)^2}$.

**Remark 6.1:** The iterations of (63) requires the recalculation of $\mathbf{F}_t$ for each time instants since $\mathbf{F}_t$ changes with time because of $\underline{\boldsymbol{\Omega}}_t$ (59). Evaluating the expectations, $\boldsymbol{\Omega}_t$ yields

$$\boldsymbol{\Omega}_t = \sqrt{\frac{2}{\pi}}\begin{bmatrix} \frac{1}{\sigma\epsilon_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\sigma\epsilon_N} \end{bmatrix} \otimes \mathbf{I}_M, \quad (64)$$

where $\sigma^2_{\epsilon_i} = E[\epsilon^2_i]$. For analytical reasons, we approximate (64) as

$$\boldsymbol{\Omega}_t \approx \sqrt{\frac{2}{\pi}}\frac{1}{(1/\sqrt{N})\sigma_{\boldsymbol{\epsilon}_t}}\mathbf{I}_{MN} \quad (65)$$

with $\sigma^2_{\boldsymbol{\epsilon}_t} = E\left[\boldsymbol{\epsilon}^T_t\boldsymbol{\epsilon}_t\right] = E\|\tilde{\psi}_t\|^2_{\boldsymbol{\xi}}$ and

$$\boldsymbol{\xi} \overset{\triangle}{=} \mathrm{bvec}\left\{\begin{bmatrix} \boldsymbol{\Lambda}_{\boldsymbol{c}} & -\boldsymbol{\Lambda}_{\boldsymbol{c}} \\ -\boldsymbol{\Lambda}_{\boldsymbol{c}} & \boldsymbol{\Lambda}_{\boldsymbol{c}} \end{bmatrix}\right\}.$$

Hence, we can calculate $\mathbf{F}_t$ by iterating the following

$$E\|\tilde{\psi}_{t+1}\|^2_{\boldsymbol{\xi}} = E\|\tilde{\psi}_t\|^2_{\boldsymbol{\xi}} - \|\tilde{\psi}_0\|^2_{\boldsymbol{\Pi}_{t-1}(\mathbf{I}-\mathbf{F}_t)\boldsymbol{\xi}} + \mathbf{b}^T\left(\mathbf{I} - \boldsymbol{\Delta}_{t-1}(\mathbf{I} - \mathbf{F}_t)\right)\boldsymbol{\xi}, \quad (66)$$

where $E\|\tilde{\psi}_0\|^2_{\boldsymbol{\xi}} = \zeta\mathbf{1}^T\boldsymbol{\Lambda}_{\boldsymbol{c}}\mathbf{1}$. In Table II, we tabulate the initial condition and the weighting matrix necessary for the recursion iterations (66) of $\sigma^2_{\boldsymbol{\epsilon}_t} = E[\boldsymbol{\epsilon}^T_t\boldsymbol{\epsilon}_t]$.

## VII. STEADY-STATE ANALYSIS

At steady-state, (36) yields

$$E\|\tilde{\psi}_\infty\|^2_{(\mathbf{I}-\mathbf{F})\boldsymbol{\sigma}} = \mathbf{b}^T\boldsymbol{\sigma}.$$

In order to calculate the steady-state performance measure $E\|\tilde{\psi}_\infty\|^2_{\boldsymbol{\sigma}'}$ we choose the weighting matrix such that

$$\boldsymbol{\sigma}' = (\mathbf{I} - \mathbf{F})\boldsymbol{\sigma}$$

then the steady-state performance measure is given by

$$E\|\tilde{\psi}_\infty\|^2_{\boldsymbol{\sigma}'} = \mathbf{b}^T(\mathbf{I} - \mathbf{F})^{-1}\boldsymbol{\sigma}'. \quad (67)$$

Similar to (67), the steady state mean square error $E[\boldsymbol{\epsilon}^T_t\boldsymbol{\epsilon}_t]$ for the single bit diffusion strategy is given by

$$E\|\tilde{\psi}_\infty\|^2_{\boldsymbol{\xi}} = \mathbf{b}^T\left(\mathbf{I} - \mathbf{F}_\infty\right)^{-1}\boldsymbol{\xi}. \quad (68)$$

We point out that $\mathbf{F}_\infty$ depends on $E\|\tilde{\psi}_\infty\|^2_{\boldsymbol{\xi}}$. Once we calculate $\mathbf{F}_\infty$ numerically by (68) or through rough approximations, we can obtain any steady state performance by (67).

## VIII. TRACKING PERFORMANCE

The diffusion implementation improves the ability of the network to track variations in the underlying statistical profiles [6]. In this section, we analyze the tracking performance of the compressive diffusion strategies in a non-stationary environment. We assume a first-order random walk model, which is commonly used in the literature [5], for $\mathbf{w}_{\mathbf{o}}(t)$ such that

$$\mathbf{w}_{\mathbf{o}}(t + 1) = \mathbf{w}_{\mathbf{o}}(t) + \mathbf{q}_t,$$

where $\mathbf{q}_t \in \mathbb{R}^M$ denotes a zero-mean vector process independent of the regression data and observation noise with covariance matrix $E[\mathbf{q}_t\mathbf{q}^T_t] = \mathbf{Q}$. We introduce the global time-variant parameter vectors as $\underline{\mathbf{w}}_{\mathbf{o}}(t) = \mathrm{col}\{\mathbf{w}_{\mathbf{o}}(t), \cdots, \mathbf{w}_{\mathbf{o}}(t)\}$ and we have the global deviation vectors as $\tilde{\phi}_t = \underline{\mathbf{w}}_{\mathbf{o}}(t) - \phi_t$ and $\tilde{\mathbf{a}}_t = \underline{\mathbf{w}}_{\mathbf{o}}(t) - \mathbf{a}_t$. Then, by (23), we obtain

$$\tilde{\psi}_{t+1} = \mathbf{X}\tilde{\psi}_t - \mathbf{DY}_t\underline{\mathbf{h}}(\mathbf{e}_t, \boldsymbol{\epsilon}_t) + \underline{\mathbf{q}}_t, \quad (69)$$

where $\underline{\mathbf{q}}_t \overset{\triangle}{=} \mathrm{col}\{\mathbf{q}_t, \cdots, \mathbf{q}_t\}$ with $2MN \times 1$ dimensions. Since we assume that $\mathbf{q}_t$ is independent from the regression data $\mathbf{u}_{i,t}$, $\mathbf{c}_{i,t}$ and the observation noise $\mathbf{v}_{i,t}$ for all $i \in \{1, \cdots, N\}$, (69) yields the following weighted-energy relation

$$E\left[\tilde{\psi}^T_{t+1}\boldsymbol{\Sigma}\tilde{\psi}_{t+1}\right] = E\left[\tilde{\psi}^T_t\mathbf{X}^T\boldsymbol{\Sigma}\mathbf{X}\tilde{\psi}_t\right]$$
$$- E\left[\tilde{\psi}^T_t\mathbf{X}^T\boldsymbol{\Sigma}\mathbf{DY}_t\underline{\mathbf{h}}(\mathbf{e}_t, \boldsymbol{\epsilon}_t)\right]$$
$$- E\left[\underline{\mathbf{h}}^T(\mathbf{e}_t, \boldsymbol{\epsilon}_t)\mathbf{Y}^T_t\mathbf{D}\boldsymbol{\Sigma}\mathbf{X}\tilde{\psi}_t\right]$$
$$+ E\left[\underline{\mathbf{h}}^T(\mathbf{e}_t, \boldsymbol{\epsilon}_t)\mathbf{Y}^T_t\mathbf{D}\boldsymbol{\Sigma}\mathbf{DY}_t\underline{\mathbf{h}}(\mathbf{e}_t, \boldsymbol{\epsilon}_t)\right]$$
$$+ E\left[\underline{\mathbf{q}}_t{}^T\boldsymbol{\Sigma}\underline{\mathbf{q}}_t\right]. \quad (70)$$

We note that (70) is similar to (40) except for the last term $E\left[\underline{\mathbf{q}}_t{}^T\boldsymbol{\Sigma}\underline{\mathbf{q}}_t\right]$. We denote $2N \times 2N$ matrix whose terms are 1 as $\underline{\mathbf{1}}_{2N} \overset{\triangle}{=} [\mathbf{1}, \cdots, \mathbf{1}]$. Then, the last term in (70) is given by $\boldsymbol{\rho}^T\boldsymbol{\sigma}$ where $\boldsymbol{\rho} = \mathrm{bvec}\{\underline{\mathbf{1}}_{2N} \otimes \mathbf{Q}\}$. Through (70), we get

$$E\|\tilde{\psi}_{t+1}\|^2_{\boldsymbol{\sigma}} = E\|\tilde{\psi}_t\|^2_{\mathbf{F}_t\boldsymbol{\sigma}} + \mathbf{b}^T\boldsymbol{\sigma} + \boldsymbol{\rho}^T\boldsymbol{\sigma}. \quad (71)$$

We define $\mathbf{F}_t$ in (37) and (59) for scalar and single-bit diffusion strategies, respectively. Similarly, $\mathbf{b}$ is introduced in (35) and (60) for the scalar (time-invariant) and single-bit diffusion strategies. We point out that (71) is different from

(36) and (58) only for the term $\boldsymbol{\rho}^T\boldsymbol{\sigma}$. As a result, at steady state, (67) and (71) leads

$$E\|\tilde{\boldsymbol{\psi}}_\infty\|_{\boldsymbol{\sigma}}^2 = (\mathbf{b} + \boldsymbol{\rho})^T(\mathbf{I} - \mathbf{F}_\infty)^{-1}\boldsymbol{\sigma}. \tag{72}$$

Through (72) and Table I, we can obtain the tracking performance of the network for the conventional performance measures. We point out that in the full diffusion configuration, $\boldsymbol{\rho} = \text{bvec}\{\underline{\mathbf{1}}_N \otimes \mathbf{Q}\}$.

In the next section, we introduce the confidence parameter and the adaptive combination method, which provides better trade-off in terms of transient and steady-state performance.

## IX. CONFIDENCE PARAMETER AND ADAPTIVE COMBINATION

The cooperation among the nodes is not beneficial in general unless the cooperation rule is chosen properly [1]. For example, uniform [18], the Metropolis [19], relative-degree rules [8] and adaptive combiners [21] provide improved convergence performance relative to the no-cooperation configuration in which nodes aim to estimate the parameter of interest $\mathbf{w_o}$ without information exchange. However, the compressive diffusion strategies have a different diffusion protocol than the full diffusion configuration. At each node $i$, we combine the local estimates $\boldsymbol{\phi}_{i,t}$ with the constructed estimates $\mathbf{a}_{j,t}$ that track the local estimates $\boldsymbol{\phi}_{j,t}$ of the neighboring nodes, i.e., $j \in \mathcal{N}_i \setminus i$. Especially at the early stages of the adaptation, the constructed estimates carry far less information than the local estimates since they are not sufficiently close to the original estimates in the mean square sense. We point out that the global deviation equation of $\boldsymbol{\phi}_t$ could be written as

$$\tilde{\boldsymbol{\phi}}_{t+1} = \left(\mathbf{I} - \mathbf{M}\mathbf{U}_t\mathbf{U}_t^T\right)\mathbf{G}\tilde{\boldsymbol{\phi}}_t - \mathbf{M}\mathbf{U}_t\mathbf{v}_t + \left(\mathbf{I} - \mathbf{M}\mathbf{U}_t\mathbf{U}_t^T\right)\mathbf{G_C}\Delta\mathbf{a}_t, \tag{73}$$

where $\Delta\mathbf{a}_t \triangleq \boldsymbol{\phi}_t - \mathbf{a}_t$. In (73), we observe that the compressive diffusion update includes one additional term, i.e., the last term on RHS of (73), different from the full diffusion configuration. We can weaken the weight of the last term by arranging the combination matrix accordingly. Hence, we add one more freedom of dimension to the update by introducing a confidence parameter $\delta$. The confidence parameter determines the weight of the local estimates relative to the constructed estimates such that the new combination matrix $\boldsymbol{\Gamma}'$ is given by

$$\boldsymbol{\Gamma}' = \delta\mathbf{I}_N + (1 - \delta)\boldsymbol{\Gamma} \tag{74}$$

where $0 \le \delta \le 1$. We note that $\delta = 1$, in which case we are confident with the local estimates, yields the no-cooperation scheme and $\delta = 0$ is the full diffusion configuration where we thrust the diffused information totally.

For the new combination matrix (74), the combination of the local estimate and the constructed estimates (11) yields

$$\mathbf{w}_{i,t+1} = (1 - \delta)\underbrace{\left[\gamma_{i,i}\boldsymbol{\phi}_{i,t+1} + \sum_{j \in \mathcal{N}_i \setminus i}\gamma_{i,j}\mathbf{a}_{j,t+1}\right]}_{\hat{\boldsymbol{\phi}}_{i,t+1}}$$
$$+ \delta\boldsymbol{\phi}_{i,t+1} \tag{75}$$

We note that (75) is a convex combination of the parameter vectors $\hat{\boldsymbol{\phi}}_{i,t+1}$ and $\boldsymbol{\phi}_{i,t+1}$. Hence, we can adapt the convex combination weight $\delta$ using a stochastic gradient update [26]–[29]. Then, (75) yields

$$\mathbf{w}_{i,t+1} = \delta_{i,t+1}\boldsymbol{\phi}_{i,t+1} + (1 - \delta_{i,t+1})\hat{\boldsymbol{\phi}}_{i,t+1}. \tag{76}$$

In [27], authors update the combination weight indirectly through a sigmoidal function. Similarly, we re-parameterize the confidence parameter $\delta_{i,t}$ using the sigmoidal function [30] and an unconstrained variable $\alpha_{i,t}$ such that

$$\delta_{i,t} = \frac{1}{1 + e^{-\alpha_{i,t}}}. \tag{77}$$

We train the unconstrained weight $\alpha_{i,t}$ using a stochastic gradient update minimizing $e_{i,t}^2 = \left(d_{i,t} - \mathbf{u}_{i,t}^T\mathbf{w}_{i,t}\right)^2$ as follows

$$\alpha_{i,t+1} = \alpha_{i,t} - \frac{1}{2}\mu_{\text{cvx}}\frac{\partial e_{i,t}^2}{\partial \alpha_{i,t}}$$
$$= \alpha_{i,t} + \mu_{\text{cvx}}e_{i,t}\mathbf{u}_{i,t}^T(\boldsymbol{\phi}_{i,t} - \hat{\boldsymbol{\phi}}_{i,t})\delta_{i,t}(1 - \delta_{i,t}). \tag{78}$$

As a result, we combine the local and constructed estimates via (76), (77) and (78).

In the next section, we provide numerical examples showing the match of the theoretical and simulated results, and the improved convergence performance with the adaptive confidence parameter.

## X. NUMERICAL EXAMPLES

In this section, we examine two distinct network scenarios where we demonstrate that the theoretical analysis accurately model the simulated results and confidence parameter provides significantly improved convergence performance. In the first example, we have a network of $5$ nodes where at each node $i$, we observe a stationary data $d_{i,t} = \mathbf{u}_{i,t}^T\mathbf{w_o} + v_{i,t}$ for $i \in \{1, 2, \cdots, N\}$. The regression data $\mathbf{u}_{i,t}$ is zero-mean Gaussian with randomly chosen standard deviation $\sigma_{u_i}$, i.e., $\sigma_{u_i} = 0.1(\sqrt{10} - 1)\mathcal{U}[0,1] + 0.1$. The variance of the observation noise is $\sigma_{n_i}^2 = 10^{-3}$. In other words, the signal-to-noise ratio over the network varies around 10 to 100. The standard deviation of the projection operator is $\sigma_{c_i} = 1$. The parameter of interest $\mathbf{w_o} \in \mathbb{R}^4$ is randomly chosen. Note that we examine a relatively small network with short filter length since the computational complexity of the theoretical performance relations (39) and (63) increases exponentially with the filter length $M$ and the network size $N$.

In the no-cooperation configuration, the combination matrix is given by $\boldsymbol{\Gamma}_0 = \mathbf{I}_N$. We use the Metropolis combination rule [19] for the full diffusion configuration where the adjacency matrix of the network is given by

$$\begin{bmatrix} 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 \end{bmatrix}.$$

In the Metropolis rule [20], the combination weights are chosen according to

$$\lambda_{i,j} = \begin{cases} \frac{1}{\max\{n_i, n_j\}} & \text{if } j \in \mathcal{N}_i \setminus i, \\ 0 & \text{if } j \notin \mathcal{N}_i, \\ 1 - \sum_{j \in \mathcal{N}_i \setminus i}\lambda_{i,j} & \text{if } i = j, \end{cases}$$

Fig. 4: Comparison of global MSD curves $1/NE\|\tilde{\phi}_t\|^2$ of the single-bit and scalar diffusion approaches for $\delta = 0$ and $\delta = 0.9$.



Fig. 5: Comparison of the global MSD curves of the no-cooperation, single-bit, scalar and full diffusion configurations in the CTA diffusion strategy.



Fig. 6: Comparison of the global EMSE curves of the no-cooperation, single-bit, scalar and full diffusion configurations in the CTA diffusion strategy.



Fig. 7: Comparison of the global MSD curves of the no-cooperation, single-bit, scalar and full diffusion configurations in the ATC diffusion strategy.



Fig. 8: The MSD curves of the construction estimate $1/NE\|\tilde{\mathbf{a}}_t\|^2$ of the single-bit and scalar diffusion approaches.

where $n_i$ and $n_j$ denote the number of neighboring nodes for $i$ and $j$. For single-bit and one-dimension diffusion strategies we examine the convergence performance for the confidence parameter $\delta = 0$ and $\delta = 0.9$ in Fig. 4. We choose the step sizes the same for the distributed LMS update (14) of all configurations at all nodes, i.e., $\mu_i = 0.042$. At each node, the step sizes for the construction update (15) are $\eta_i = 0.0015$ (for single-bit approach) and $\eta_i = 0.25$ (for one-dimension diffusion approach). For the single-bit diffusion approach, we set $\zeta = 0.001$ to initialize $\mathbf{a}_{j,t}$.

In Fig. 4, we show the global MSD curves, i.e., $E\|\tilde{\phi}_t\|^2$, of the single-bit and scalar diffusion approaches and compare the performance for different $\delta$ values. The confidence parameter $\delta = 0.9$ implies that we give ten times more weight to the local estimate $\phi_{i,t}$ than the constructed estimates $\mathbf{a}_{j,t}$ where $j \in \mathcal{N}_i \setminus i$. The Fig. 4 demonstrates that the confidence parameter $\delta = 0.9$ improves the convergence performance of the compressive diffusion strategies. For the same example, Fig. 5, Fig. 6 and Fig. 7 compare the convergence performance of single-bit and scalar diffusion strategies with the no-cooperation and full diffusion configurations for $\delta = 0.9$, which shows the match of the theoretical and ensemble averaged (we perform 200 independent trials) performance results. The Fig. 5 and Fig. 6 show the time-evolution of the MSD

and EMSE curves in the CTA diffusion strategy while the Fig. 7 displays the time-evolution of the MSD curves in the ATC diffusion strategy in which the theoretical curves (39) and (63) are iterated according to the Table I and II. We note that we obtain similar MSD curves in the CTA and ATC strategies. Since we set $\delta = 0.9$ and the outcomes of the adaptation and combination operations contain relatively close amount of information.

The Fig. 8 demonstrates the convergence of the constructed

Fig. 9: Network topology with $N = 20$ for the Example 2.



Fig. 10: Network statistical profile ($\sigma_{n_i} = 0.1$).



Fig. 11: The global EMSE curves in relatively large network size and long filter length and the confidence parameter is adaptively chosen for single-bit and scalar diffusion strategies.

estimates $\mathbf{a}_{j,t}$'s to the parameter of interest $\mathbf{w_o}$ in the mean-square sense. We point out that the recursions (39) and (63) also provide the global mean-square deviation of the constructed estimates for the certain combination weight $\boldsymbol{\Sigma}$ in Table I and the theoretical recursion matches with the simulated results.

In the second example, we examine the convergence performance of the adaptive confidence parameter in relatively large network $N = 20$ with long filter length $M = 10$ (See Fig. 9). We again observe a stationary data $d_{i,t} = \mathbf{u}_{i,t}^T \mathbf{w_o} + v_{i,t}$ for $i \in \{1, 2, \cdots, N\}$. The regressor data $\mathbf{u}_{i,t}$ is zero-mean i.i.d. Gaussian whose standard deviation varies over the network as in Fig. 10. The observation noise $v_{i,t}$ is zero-mean i.i.d. Gaussian whose variance is $\sigma_{n_i} = 10^{-2}$. We note that the signal-to-noise ratio varies from 10 to 100 over the network similar to the example 1. The standard deviation of the projection operator $\mathbf{c}_{i,t}$ is $\sigma_{c_i} = 1$ and the parameter of interest $\mathbf{w_o} \in \mathbb{R}^{10}$ is randomly chosen.

We again use the Metropolis rule as the combination rule, however, in this example, we adapt the confidence parameter through (77) and (78) where we resort to the convex mixture of the adaptive filtering algorithms [26]–[29]. We also choose the step sizes the same for the distributed LMS update (14) of all configurations at all nodes, i.e., $\mu_i = 0.042$. In example 2, the

step sizes for the construction update (15) are $\eta_i = 0.0042$ (for single-bit approach) and $\eta_i = 0.1$ (for one-dimension diffusion approach). We set $\mu_{\mathrm{cvx}} = 10$ in (78). The Fig. 11 shows the global MSD curves of the no-cooperation, single-bit, scalar and full diffusion strategies. We observe that the adaptive confidence parameter improves the convergence performance of the compressive diffusion strategies far more such that they achieve comparable performance while the reduction of the communication load is tremendous.

## XI. CONCLUSION

In the diffusion based distributed estimation strategies, the communication load increases far more in the large networks or highly connected network of nodes. Hence, the compressive diffusion approach plays an essential role in achieving comparable convergence performance with the full diffusion configurations while reducing the communication load significantly. We provide a complete performance analysis for the compressive diffusion strategies. We analyze the mean-square convergence, steady-state behavior and the tracking performance of the scalar and single-bit diffusion approaches. The numerical examples show the theoretical analysis model the simulated results accurately. Additionally, we introduce the confidence parameter concept, which adds one more freedom of dimension to the combination rule in order to improve the convergence performance. When we adapt the confidence parameter using the well-known mixture algorithms, we observe enormous enhancement in the convergence performance of the compressive diffusion strategies even for the relatively long filter lengths.

## APPENDIX A
## PROOF FOR LEMMA 1

We first show the equality of (43) for the two-node case. Then the extension for a larger network is straight forward. We can rewrite the term on the left hand side (LHS) of (43) as

$$E[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X_u}^T \boldsymbol{\Sigma}_2 \mathbf{N C}_t \mathrm{sign}(\boldsymbol{\epsilon}_t)]$$

$$= E\left[ \tilde{\boldsymbol{\psi}}_t^T \mathbf{X_u}^T \underbrace{\begin{bmatrix} \boldsymbol{\varsigma_1} & \boldsymbol{\varsigma_2} \\ \boldsymbol{\varsigma_3} & \boldsymbol{\varsigma_4} \end{bmatrix}}_{\boldsymbol{\Sigma}_2} \mathbf{N C}_t \mathrm{sign}(\boldsymbol{\epsilon}_t) \right]. \quad (79)$$

After some algebra, (79) yields

$$\begin{aligned}
E[\tilde{\boldsymbol{\psi}}_t^T & \mathbf{X_u}^T \boldsymbol{\Sigma}_2 \mathbf{N C}_t \mathrm{sign}(\boldsymbol{\epsilon}_t)] \\
= & E[(\gamma_{11} \tilde{\boldsymbol{\phi}}_{1,t}^T + \gamma_{12} \tilde{\mathbf{a}}_{2,t}^T) \boldsymbol{\varsigma_1} \eta_1 \mathbf{c}_{1,t} \mathrm{sign}(\epsilon_{1,t})] \\
& + E[(\gamma_{11} \tilde{\boldsymbol{\phi}}_{1,t}^T + \gamma_{12} \tilde{\mathbf{a}}_{2,t}^T) \boldsymbol{\varsigma_2} \eta_2 \mathbf{c}_{2,t} \mathrm{sign}(\epsilon_{2,t})] \\
& + E[(\gamma_{22} \tilde{\boldsymbol{\phi}}_{2,t}^T + \gamma_{21} \tilde{\mathbf{a}}_{1,t}^T) \boldsymbol{\varsigma_3} \eta_1 \mathbf{c}_{1,t} \mathrm{sign}(\epsilon_{1,t})] \\
& + E[(\gamma_{22} \tilde{\boldsymbol{\phi}}_{2,t}^T + \gamma_{21} \tilde{\mathbf{a}}_{1,t}^T) \boldsymbol{\varsigma_4} \eta_2 \mathbf{c}_{2,t} \mathrm{sign}(\epsilon_{2,t})].
\end{aligned}$$
$$(80)$$

In order to evaluate the expectations on the RHS of (80), we assume that the step sizes are sufficiently small and filter is sufficiently long so that the deviation terms changes negligibly

slow with respect to the regressor data $\mathbf{c}_{i,t}$. Then, according to the Price's result [31], [32], we obtain

$$E[\tilde{\boldsymbol{\psi}}_t^T \mathbf{X_u}^T \boldsymbol{\Sigma}_2 \mathbf{N} \mathbf{C}_t \mathrm{sign}(\boldsymbol{\epsilon}_t)]$$

$$= E[(\gamma_{11}\tilde{\boldsymbol{\phi}}_{1,t}^T + \gamma_{12}\tilde{\mathbf{a}}_{2,t}^T)\boldsymbol{\varsigma_1}\eta_1\mathbf{c}_{1,t}\epsilon_{1,t}]\frac{E|\epsilon_{1,t}|}{E[\epsilon_{1,t}^2]}$$

$$+ E[(\gamma_{11}\tilde{\boldsymbol{\phi}}_{1,t}^T + \gamma_{12}\tilde{\mathbf{a}}_{2,t}^T)\boldsymbol{\varsigma_2}\eta_2\mathbf{c}_{2,t}\epsilon_{2,t}]\frac{E|\epsilon_{2,t}|}{E[\epsilon_{2,t}^2]}$$

$$+ E[(\gamma_{22}\tilde{\boldsymbol{\phi}}_{2,t}^T + \gamma_{21}\tilde{\mathbf{a}}_{1,t}^T)\boldsymbol{\varsigma_3}\eta_1\mathbf{c}_{1,t}\epsilon_{1,t}]\frac{E|\epsilon_{1,t}|}{E[\epsilon_{1,t}^2]}$$

$$+ E[(\gamma_{22}\tilde{\boldsymbol{\phi}}_{2,t}^T + \gamma_{21}\tilde{\mathbf{a}}_{1,t}^T)\boldsymbol{\varsigma_4}\eta_2\mathbf{c}_{2,t}\epsilon_{2,t}]\frac{E|\epsilon_{2,t}|}{E[\epsilon_{2,t}^2]}. \quad (81)$$

Rearranging (81) into a matrix product form leads (43). Following the same way, we can also get (44) and the proof is concluded.

## APPENDIX B
## PROOF FOR LEMMA 2

We derive the RHS of (49) for the two-node case for simplicity, however, it also satisfies any order of network. For two-node case, the LHS of (49) yields

$$E\left[\mathrm{sign}(\boldsymbol{\epsilon}_t)^T \mathbf{C}_t^T \mathbf{N}\boldsymbol{\Sigma}_4 \mathbf{N}\mathbf{C}_t\mathrm{sign}(\boldsymbol{\epsilon}_t)\right]$$

$$= E\left[\mathrm{sign}(\epsilon_{1,t})\mathbf{c}_{1,t}^T\eta_1\boldsymbol{\varsigma_1}\eta_1\mathbf{c}_{1,t}\mathrm{sign}(\epsilon_{1,t})\right]$$

$$+ E\left[\mathrm{sign}(\epsilon_{1,t})\mathbf{c}_{1,t}^T\eta_1\boldsymbol{\varsigma_2}\eta_2\mathbf{c}_{2,t}\mathrm{sign}(\epsilon_{2,t})\right]$$

$$+ E\left[\mathrm{sign}(\epsilon_{2,t})\mathbf{c}_{2,t}^T\eta_2\boldsymbol{\varsigma_3}\eta_1\mathbf{c}_{1,t}\mathrm{sign}(\epsilon_{1,t})\right]$$

$$+ E\left[\mathrm{sign}(\epsilon_{2,t})\mathbf{c}_{2,t}^T\eta_2\boldsymbol{\varsigma_4}\eta_2\mathbf{c}_{2,t}\mathrm{sign}(\epsilon_{2,t})\right].$$

We re-emphasize that the regressor $\mathbf{c}_{i,t}$ is spatially and temporarily independent. Hence, we obtain

$$E\left[\mathrm{sign}(\boldsymbol{\epsilon}_t)^T \mathbf{C}_t^T \mathbf{N}\boldsymbol{\Sigma}_4 \mathbf{N}\mathbf{C}_t\mathrm{sign}(\boldsymbol{\epsilon}_t)\right]$$

$$= E\left[\mathbf{c}_{1,t}^T\eta_1\boldsymbol{\varsigma_1}\eta_1\mathbf{c}_{1,t}\right] + E\left[\mathbf{c}_{2,t}^T\eta_2\boldsymbol{\varsigma_4}\eta_2\mathbf{c}_{2,t}\right]$$

$$+ E\left[\mathbf{c}_{1,t}\mathrm{sign}(\epsilon_{1,t})\right]^T\eta_1\boldsymbol{\varsigma_2}\eta_2 E\left[\mathbf{c}_{2,t}\mathrm{sign}(\epsilon_{2,t})\right]$$

$$+ E\left[\mathbf{c}_{2,t}\mathrm{sign}(\epsilon_{2,t})\right]^T\eta_2\boldsymbol{\varsigma_3}\eta_1 E\left[\mathbf{c}_{1,t}\mathrm{sign}(\epsilon_{1,t})\right]. \quad (82)$$

Using Price's result, we can evaluate the last two terms on the RHS of (82) as follows

$$E\left[\mathbf{c}_{1,t}\mathrm{sign}(\epsilon_{1,t})\right] = \frac{E|\epsilon_{1,t}|}{E[\epsilon_{1,t}^2]}E\left[\mathbf{c}_{1,t}\epsilon_{1,t}\right]$$

and

$$E\left[\mathbf{c}_{2,t}\mathrm{sign}(\epsilon_{2,t})\right] = \frac{E|\epsilon_{2,t}|}{E[\epsilon_{2,t}^2]}E\left[\mathbf{c}_{2,t}\epsilon_{2,t}\right].$$

We point out that the terms involving the diagonal entries of the weighting matrix $\boldsymbol{\Sigma}_4$ in (82) does not include the deviation terms. As a result, rearranging (82) into a compact form results in (49). This concludes the proof.

## REFERENCES

[1] A. H. Sayed, S.-Y. Tu, J. Chen, X. Zhao, and Z. J. Towfic, "Diffusion strategies for adaptation and learning over networks: an examination of distributed strategies and network behavior," *IEEE Signal Processing Magazine*, vol. 30, no. 3, pp. 155–171, 2013.

[2] D. Li, K. D. Wong, Y. H. Hu, and A. M. Sayeed, "Detection, classification, and tracking of targets," *IEEE Signal Processing Magazine*, vol. 19, no. 2, pp. 17–29, 2002.

[3] I. Akyildiz, W. Su, Y. Sankarasubramaniam, and E. Cayirci, "A survey on sensor networks," *IEEE Communications Magazine*, vol. 40, no. 8, pp. 102–114, 2002.

[4] D. Estrin, L. Girod, G. Pottie, and M. Srivastava, "Instrumenting the world with wireless sensor networks," in *Proc. Int. Conf. Acoust., Speech, and Signal Process. (ICASSP)*, vol. 4, 2001, pp. 2033–2036 vol.4.

[5] A. H. Sayed, *Fundamentals of Adaptive Filtering*. New York: Wiley, 2003.

[6] C. G. Lopes and A. H. Sayed, "Diffusion least-mean squares over adaptive networks: Formulation and performance analysis," *IEEE Transactions on Signal Processing*, vol. 56, no. 7, pp. 3122–3136, 2008.

[7] F. S. Cattivelli and A. H. Sayed, "Diffusion LMS strategies for distributed estimation," *IEEE Transactions on Signal Processing*, vol. 58, no. 3, pp. 1035–1048, 2010.

[8] F. S. Cattivelli, C. G. Lopes, and A. H. Sayed, "Diffusion recursive least-squares for distributed estimation over adaptive networks," *IEEE Transactions on Signal Processing*, vol. 56, no. 5, pp. 1865–1877, 2008.

[9] F. S. Cattivelli and A. H. Sayed, "Diffusion strategies for distributed kalman filtering and smoothing," *IEEE Transactions on Automatic Control*, vol. 55, no. 9, pp. 2069–2084, 2010.

[10] S.-Y. Tu and A. H. Sayed, "Diffusion strategies outperform consensus strategies for distributed estimation over adaptive networks," *Signal Processing, IEEE Transactions on*, vol. 60, no. 12, pp. 6217–6234, 2012.

[11] X. Zhao and A. H. Sayed, "Performance limits for distributed estimation over lms adaptive networks," *IEEE Transactions on Signal Processing*, vol. 60, no. 10, pp. 5107–5124, 2012.

[12] M. O. Sayin and S. S. Kozat, "Single bit and reduced dimension diffusion strategies over distributed networks," *IEEE Signal Processing Letters*, vol. 20, no. 10, pp. 976–979, 2013.

[13] D. L. Donoho, "Compressed sensing," *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1289–1306, 2006.

[14] R. G. Baraniuk, V. Cevher, and M. B. Wakin, "Low-dimensional models for dimensionality reduction and signal recovery: A geometric perspective," *Proceedings of the IEEE*, vol. 98, no. 6, pp. 959–971, 2010.

[15] S. Xie and H. Li, "Distributed LMS estimation over networks with quantised communications," *International Journal of Control*, vol. 86, no. 3, pp. 478–492, 2013.

[16] A. Ribeiro, G. B. Giannakis, and S. I. Roumeliotis, "SOI-KF: Distributed kalman filtering with low-cost communications using the sign of innovations," *IEEE Transactions on Signal Processing*, vol. 54, no. 12, pp. 4782–4795, 2006.

[17] H. Sayyadi and M. R. Doostmohammadian, "Finite-time consensus in directed switching network topologies and time-delayed communications," *Scientia Iranica*, vol. 18, no. 1, pp. 75–85, February 2011.

[18] V. D. Blondel, J. M. Hendrickx, A. Olshevsky, and J. N. Tsitsiklis, "Convergence in multiagent coordination, consensus and flocking," in *Proc. Joint 44th IEEE Conf. Decision Control Eur. Control Conf. (CDC-ECC), Seville, Spain, Dec.*, 2005, pp. 2996–3000.

[19] L. Xiao and S. Boyd, "Fast linear iterations for distributed averaging," *Syst. Control Lett.*, vol. 53, no. 1, pp. 65–78, 2004.

[20] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller, "Equation of state calculations by fast computing machines," *The Journal of Chemical Physics*, vol. 21, no. 6, 1953.

[21] N. Takahashi, I. Yamada, and A. H. Sayed, "Diffusion least-mean squares with adaptive combiners: Formulation and performance analysis," *IEEE Transactions on Signal Processing*, vol. 58, no. 9, pp. 4795–4810, 2010.

[22] A. H. Sayed and C. G. Lopes, "Distributed adaptive learning mechanisms," in *Handbook on Array Processing and Sensor Networks*, S. S. Haykin and K. J. R. Liu, Eds. New York: Wiley, 2009.

[23] E. Barker and J. Kelsey, "Recommendation for random number generation using deterministic random bit generators," *NIST SP800-90A*, 2012. [Online]. Available: http://csrc.nist.gov/publications/nistpubs/800-90A/SP800-90A.pdf

[24] J. Joutsensalo and T. Ristaniemi, "Synchronization by pilot signal," in *Proc. Int. Conf. Acoust., Speech, and Signal Process. (ICASSP)*, 1999, pp. 2663–2666.

[25] T. Y. Al-Naffouri and A. H. Sayed, "Transient analysis of data-normalized adaptive filters," *IEEE Transactions on Signal Processing*, vol. 51, no. 3, pp. 639–652, 2003.

[26] J. Arenas-Garcia, V. Gomez-Verdejo, and A. R. Figueiras-Vidal, "New algorithms for improved adaptive convex combination of lms transversal filters," *IEEE Transactions on Instrumentation and Measurement*, vol. 54, no. 6, pp. 2239–2249, 2005.

[27] J. Arenas-Garcia, A. R. Figueiras-Vidal, and A. H. Sayed, "Mean-square performance of a convex combination of two adaptive filters," *IEEE Transactions on Signal Processing*, vol. 54, no. 3, pp. 1078–1090, 2006.

[28] M. T. M. Silva and V. H. Nascimento, "Improving the tracking capability of adaptive filters via convex combination," *IEEE Transactions on Signal Processing*, vol. 56, no. 7, pp. 3137–3149, 2008.

[29] S. S. Kozat, A. T. Erdogan, A. C. Singer, and A. H. Sayed, "Steady state MSE performance analysis of mixture approaches to adaptive filtering," *IEEE Transactions on Signal Processing*, vol. 58, no. 8, pp. 4050–4063, August 2010.

[30] J. Han and C. Moraga, "The influence of the sigmoid function parameters on the speed of backpropagation learning," in *From Natural to Artificial Neural Computation*, J. Mira and F. Sandoval, Eds. Springer Berlin Heidelberg, 1995.

[31] R. Price, "A useful theorem for nonlinear devices having gaussian inputs," *IEEE Transactions on Information Theory*, vol. 4, no. 2, pp. 69–72, 1958.

[32] E. McMahon, "An extension of price's theorem (corresp.)," *IEEE Transactions on Information Theory*, vol. 10, no. 2, pp. 168–168, 1964.