# Fast and Accurate Amplitude Demodulation of Wideband Signals

## Mantas Gabrielaitis

*Abstract*—**Amplitude demodulation is a classical operation used in signal processing. For a long time, its effective applications in practice have been limited to narrowband signals. In this work, we generalize amplitude demodulation to wideband signals. We pose demodulation as a recovery problem of an oversampled corrupted signal and introduce special iterative schemes belonging to the family of alternating projection algorithms to solve it. Sensibly chosen structural assumptions on the demodulation outputs allow us to reveal the high inferential accuracy of the method over a rich set of relevant signals. This new approach surpasses current state-of-the-art demodulation techniques apt to wideband signals in computational efficiency by up to many orders of magnitude with no sacrifice in quality. Such performance opens the door for applications of the amplitude demodulation procedure in new contexts. In particular, the new method makes online and large-scale offline data processing feasible, including the calculation of modulator-carrier pairs in higher dimensions and poor sampling conditions, independent of the signal bandwidth. We illustrate the utility and specifics of applications of the new method in practice by using natural speech and synthetic signals.**

*Index Terms*—**Alternating projections, amplitude demodulation, convex programming, fast algorithms, multidimensional signals, nonuniform sampling, speech processing, wideband signals.**

## I. INTRODUCTION

**A**MPLITUDE demodulation refers to the decomposition of a signal into a product of a slow-varying modulator-envelope and a fast-varying carrier. First introduced in radio communications [1], this procedure has found applications in data acquisition and processing related to a broad range of phenomena. Automatic speech recognition [2], atomic force microscopy [3], ultrasound imaging [4], brainwave [5], seismic trace [6], and fingerprint [7] analyses are a few among many examples to mention.

Originally, amplitude demodulation was intended for use with signals built of locally sinusoidal, i.e., narrowband, carriers. Several classical approaches excel in this setting, with Gabor's analytic-signal (AS) method being a long-standing champion [8], [9]. Nonetheless, many relevant problems inevitably require demodulating signals that feature wideband carriers, typically of (quasi)-harmonic, (quasi)-random, or spike-train origin [10]–[18] (see Suppl. Mat. H for an overview). When applied to them, the classical techniques fail, misleadingly mixing the carrier and modulator information [19], [20].

For a long time, no consistent and accurate way of demodulating wideband signals was known. Typically, a proxy of the modulator would be obtained by rectifying and then low-pass filtering the signal. Different implementations of

M. Gabrielaitis is with the Institute of Science and Technology Austria, 3400 Klosterneuburg, Austria (e-mail: mantas.gabrielaitis@ist.ac.at)

this procedure, each adapted for a specific signal class, were suggested (see, e.g., [17], [21], [22]). The estimates of signal modulators obtained in this way, however, are neither accurate nor consistent between different methods. The carriers and modulators are not appropriately separated either, i.e., they can be demodulated further by iterating the same procedure [23]. Moreover, the carrier estimates are often unbounded, even in well-defined situations (see, e.g., [23, Fig. 3.1]).

Recently, two promising demodulation approaches suitable to signals with arbitrary bandwidths have been formulated. Turner and Sahani shaped demodulation into a statistical inference problem [23], [24]. In this so-called probabilistic amplitude demodulation (PAD) approach, the modulator and carrier are inferred from the signal as latent variables of an appropriately selected statistical model. Mathematically, PAD defines a maximization of a posteriori probability, a high-dimensional nonlinear optimization task. In another work, Sell and Slaney chose a deterministic route to demodulation [19]. In their linear-domain convex (LDC) approach, the modulator is described as a minimum-power signal with penalized high-frequency terms lying above the original waveform. This problem is convex and thus amenable to more efficient optimization methods than the PAD.

The PAD and LDC techniques separate the modulator and carrier information of various synthetic wideband signals with a high degree of accuracy [19], [23]. The principal weakness of these approaches is a huge associated computational burden, which impedes their use in practical situations (see Section IV for the performance evaluations). In particular, online or large-scale offline signal processing is out of reach for the PAD and LDC demodulations. Besides, derivations of these methods are guided more by high-level modulator or carrier properties and computational tractability rather than strict recovery conditions. Hence, the boundaries of their validity in the context of real-world signals are somewhat blurred.

In this work, we frame demodulation as a problem of modulator recovery from an unlabeled mix of its true and corrupted sample points. We show that, under some loose constraints on carriers and modulators, high-accuracy demodulation can be achieved through exact or approximate norm minimization. We introduce different versions of custom-made alternating projection algorithms and test them in numerical experiments to solve this task. The new approach is shown to be free of the performance limitations inherent to the PAD and LDC methods. In particular, it combines the computational economy of the classical AS technique with the capacity to recover a wide range of arbitrary-bandwidth signals. We reveal the power of the new approach in terms of efficiency, accuracy, consistency, and robustness to corrupted data through theoretical analysis

and illustrate it using synthetic signals with known structure. The use of the new method in realistic online and offline settings is demonstrated by applying it to natural speech.

## II. MATHEMATICAL FORMULATION OF THE PROBLEM

In what follows, we assume the representation of a real-valued signal $s(t)$ formed by a finite collection of its values uniformly sampled over a limited time interval: $s_i \equiv s(t_i)$, $i \in \mathcal{I}_n = \{1, 2, \ldots, n\}$. Thus, a realization of the signal, $\mathbf{s} \equiv (s_1, s_2, \ldots, s_n)^T$, is an element of an $n$-dimensional Euclidean space $\mathbb{R}^n$, i.e., a linear space equipped with the inner product $\langle \mathbf{s}^{(1)}, \mathbf{s}^{(2)} \rangle = \sum_{i=1}^{n} (s_i^{(1)} \cdot s_i^{(2)})$, which induces the Euclidean norm $\|\mathbf{s}\|_2 = \sqrt{\langle \mathbf{s}, \mathbf{s} \rangle}$. We use modulo $n$ arithmetic for indexes of vector components in this work.

### A. Demodulation constraints

The task of demodulation is to factorize a signal $\mathbf{s} \in \mathbb{R}^n$ into a modulator $\mathbf{m} \in \mathbb{R}^n$ and a carrier $\mathbf{c} \in \mathbb{R}^n$:

$$\mathbf{s} = \mathbf{m} \circ \mathbf{c}, \tag{1}$$

where symbol $\circ$ denotes an elementwise product of two vectors. There exists an uncountable number of pairs of $\mathbf{m}$ and $\mathbf{c}$ that satisfy (1). Thus, further constraints are needed to define its unique solution. It is precisely these constraints that give a distinct character to different demodulation methods and set the domain of their validity [9], [19], [23], [25].

In this work, we introduce the extra demodulation restrictions by imposing some general assumptions on $\mathbf{m}$ and $\mathbf{c}$.

We define feasible modulators as elements of a convex set

$$\mathcal{M}_\omega = \mathcal{S}_{\geq \mathbf{0}} \cap \mathcal{S}_\omega, \tag{2}$$

where

$$\begin{aligned}
\mathcal{S}_{\geq \mathbf{0}} &= \{\mathbf{x} \in \mathbb{R}^n : x_i \geq 0, \ i \in \mathcal{I}_n\}, \\
\mathcal{S}_\omega &= \{\mathbf{x} \in \mathbb{R}^n : (\mathbf{Fx})_i = 0, \ i \in (\mathcal{I}_n \setminus \mathcal{I}_n^\omega)\}, \\
\mathcal{I}_n^\omega &= \{i \in \mathcal{I}_n : i \leq \omega\} \cup \{i \in \mathcal{I}_n : i > n + 1 - \omega\}.
\end{aligned} \tag{3}$$

In (3), $\mathbf{F}$ denotes the operator of the unitary discrete Fourier transform (DFT), and $(\ldots)_i$ marks the $i$-th component of the argument vector. Hence, in our framework, modulators are nonnegative low-pass signals whose rate of variation is limited by the cutoff frequency $\omega$ (with $1 \leq \omega \leq \lceil n/2 \rceil$), which parametrizes $\mathcal{M}_\omega$. This is a formal definition of the classical modulator-envelope [1], [26].

We declare feasible carriers as elements of a nonconvex set

$$\mathcal{C}_d = \mathcal{S}_{|\cdot| \leq \mathbf{1}} \cap \mathcal{S}_{\{1\}, d}, \tag{4}$$

where

$$\begin{aligned}
\mathcal{S}_{|\cdot| \leq \mathbf{1}} &= \{\mathbf{x} \in \mathbb{R}^n : |x_i| \leq 1, \ i \in \mathcal{I}_n\}, \\
\mathcal{S}_{\{1\}, d} &= \Big\{\mathbf{x} \in \mathbb{R}^n : (\forall i) \sum_{j=i}^{i+d-1} I_{\{1\}}(|x_j|) \geq 1, \\
&\qquad (\exists i) \sum_{j=i}^{i+d-1} I_{\{1\}}(|x_j|) = 1\Big\},
\end{aligned} \tag{5}$$

with $I_{\{1\}}$ being the indicator function of the singleton $\{1\}$. The set $\mathcal{S}_{|\cdot| \leq \mathbf{1}}$ implies the boundedness of $\mathbf{c}$ between $-1$ and $1$. This restriction follows from the standard notion that the time-dependent amplitude of an amplitude-modulated $\mathbf{s}$ is purely set

by $\mathbf{m}$. Meanwhile, $\mathcal{S}_{\{1\}, d}$ fixes to $d$ the maximum gap between any two neighboring components of $\mathbf{c}$ whose absolute values are equal to $1$.[1] As shown next, this constraint allows formulating extensive demodulation guarantees while only moderately affecting the scope of relevant carriers. Bandwidth-wise, $\mathcal{C}_d$ covers the whole range, from zero (sinusoidal) to flat (random spike) bandwidth signals, and defines the qualifier "wideband" used in this work. Note that the bandwidth of $\mathbf{c} \in \mathcal{C}_d$ is mostly determined not by $d$ but by the arrangement of the $|c_i| = 1$ and other sample points.[2] Instead, as we see next, $d$ decides whether a chosen $\mathbf{c} \in \mathcal{C}_d$ can be restored after modulation.

### B. Demodulation as modulator recovery

Note that, assuming $\mathbf{c} \in \mathcal{C}_d$, $|\mathbf{s}|$ can be seen as a corrupted version of $\mathbf{m}$: $|s_i| = m_i$ when $|c_i| = 1$, and $|s_i| \neq m_i$ otherwise. Further, if $\mathbf{m}$ can be found from $|\mathbf{s}|$, $\mathbf{c}$ follows from (1) uniquely ($c_i = s_i/m_i$), except the sample points with $m_i = 0$. The latter, if any, are sparse and can be typically interpolated from the neighboring points. Hence, in our case, demodulation is virtually a problem of reconstructing $\mathbf{m}$ from a mix of its *true* ($i : |s_i| = m_i$) and *corrupted* ($i : |s_i| \neq m_i$) sample points when the class of each point is unknown. This viewpoint is at the core of the developments that follow next.

### C. Modulator recovery through norm minimization

Our approach to demodulation builds around the estimator

$$\hat{\mathbf{m}} = \underset{\mathbf{x} \in \mathcal{S}_{\geq |\mathbf{s}|} \cap \mathcal{S}_\varpi}{\arg \min} \|\mathbf{x}\|_2, \tag{6}$$

where $\mathcal{S}_{\geq |\mathbf{s}|} = \{\mathbf{x} \in \mathbb{R}^n : x_i \geq |s_i|, i \in \mathcal{I}_n\}$. Note that $\mathbf{m} \in \mathcal{S}_{\geq |\mathbf{s}|} \cap \mathcal{S}_\varpi$ if $\varpi \geq \omega$. The restriction corresponding to $\mathcal{S}_{\geq |\mathbf{s}|}$ assures that $\mathbf{x}$ does not fall below $\mathbf{m}$ at the true sample points, i.e., points where $m_i = |s_i|$. If, besides, the true sample points are spread densely enough, we expect the norm minimization to enforce $\hat{m}_i = m_i$ at these points. But then, $\hat{\mathbf{m}} = \mathbf{m}$ by the discrete sampling theorem. The foundation for this intuitive consideration is laid by the following results (see Suppl. Mat. B for the proofs).

**Proposition II.1.** *For almost every* $\mathbf{m} \in \mathcal{M}_\omega$, $\hat{\mathbf{m}} = \mathbf{m}$ *only if* $\varpi \geq \omega$, *and* $\mathbf{c} \in \mathcal{C}_d$ *with* $n_s \equiv \sum_{i=1}^{n} I_{\{1\}}(|c_i|) \geq \varpi + \omega - 1 \implies d \leq n - (\varpi + \omega - 2)$.[3]

**Proposition II.2.** *Consider* $\mathbf{m} \in \mathcal{M}_\omega$ *and* $\tilde{\mathbf{c}} \in \mathcal{C}_{\tilde{d}}$ *with* $|\tilde{c}_i| = 1$ *for* $i \in \mathcal{J}_n \subseteq \mathcal{I}_n$, *and* $\tilde{c}_i = 0$ *otherwise. If* $\hat{\mathbf{m}} = \mathbf{m}$ *holds for the* $\mathbf{m}$ *and* $\tilde{\mathbf{c}}$, *then it also holds for every pair made of the same* $\mathbf{m}$ *and any* $\mathbf{c} \in \mathcal{C}_d$ *with* $d \leq \tilde{d}$ *and* $|c_i| = 1$ *for* $i \in \mathcal{J}_n$.

**Proposition II.3.** *Assume* $\mathbf{m} \in \mathcal{M}_\omega$ *and* $\mathbf{c} \in \mathcal{C}_d$ *with* $\varpi \geq \omega$. *If, additionally, there exist* $d \in \mathcal{I}_n$ *and* $i \in \mathcal{I}_d$ *such that* $n_s \equiv$

---

[1] The requirement of the existence of at least one gap of length $d$ in the definition of $\mathcal{S}_{\{1\}, d}$ assures that $\mathcal{C}_{d_1} \cap \mathcal{C}_{d_2} = \emptyset$ if $d_1 \neq d_2$. Such parametrization of the carrier set allows specifying more definite demodulation conditions.

[2] For example, even $\mathcal{C}_1$, which features the most limited repertoire among all $\mathcal{C}_d$, has zero-bandwidth elements (consider the $\mathbf{c}$ with $c_i = (-1)^i$) and elements with approximately flat amplitude spectra (consider a $\mathbf{c}$ with $c_i$ randomly chosen from $\{-1, 1\}$).

[3] In fact, as follows from the proof of this proposition in Suppl. Mat. B, the condition that $\mathbf{c} \in \mathcal{C}_d$ for at least some $d$ is necessary for strictly every $\mathbf{m}$.

$(n/d) \in \mathbb{N}_+$, $n_s \geq \varpi + \omega - 1$, and $|c_{i+(j-1)\cdot d}| = 1$ for every $j \in \mathcal{I}_{n_s}$, then $\hat{\mathbf{m}} = \mathbf{m}$.

*Proposition II.1* reveals the tight match of $\hat{\mathbf{m}}$ to $\mathcal{C}_d$: no $\mathbf{m} \in \mathcal{M}_\omega$ can be inferred from $\mathbf{s}$ by $\hat{\mathbf{m}}$ precisely if $\mathbf{c} \notin \mathcal{C}_d$. It also establishes the central role of the presence of true sample points in the recovery: for almost every $\mathbf{m} \in \mathcal{M}_\omega$, at least the number $\varpi + \omega - 1$ of such points is needed. *Proposition II.2* further consolidates the latter view by stating that the success of the exact recovery of an $\mathbf{m} \in \mathcal{M}_\omega$ via $\hat{\mathbf{m}}$ is fully determined by the number and positions of the true sample points. In particular, if exact demodulation is possible for some $\tilde{\mathbf{c}}$ with $\tilde{c}_i \in \{0, 1\}$, then it is possible for any $\mathbf{c}$ with $|c_i| = 1$ at $i \in \{j : \tilde{c}_j = 1\}$ independent of other sample points.

In *Proposition II.1*, $\varpi \geq \omega$ and $n_s \geq \varpi + \omega - 1$ imply $n_s \geq 2\omega - 1$, which is a sufficient condition for $\mathbf{m}$ recovery in the classical setup when all true sample points are known (see the remark below the proof of *Proposition A.1* in Suppl. Mat. A). Hence, the data corruption manifesting in our problem necessitates further constraints on the number or positions of true sample points. In particular, *Proposition II.3* certifies a full recovery of $\mathbf{m}$ if $\varpi \geq \omega$, and there exists a (not necessarily known) subset of at least $\varpi + \omega - 1$ regularly-spaced true sample points. The latter condition covers a wide range of practically relevant carriers, including: (1) the classical $\sin(2\pi\nu\mathbf{t} + \phi)$ with $\nu \geq \omega$, (2) harmonic signals, (3) regular spike-trains of $|c_i| = 1$. More generally, any (non)stationary time-series with regularly placed $|c_i| = 1$ regardless of the remaining points are eligible.

In addition to the regularity of true sample points, *Proposition II.3* requires $n/d$ to be an integer. Nevertheless, numerical experiments reveal that both of these conditions can be ignored without practically relevant consequences (see Suppl. Mat. C and Fig. 10 there). In particular, we found that the discrepancy between $\mathbf{m}$ and $\hat{\mathbf{m}}$ is vanishing with an overwhelming probability for any $\mathbf{c} \in \mathcal{C}_d$ if $\varpi \geq \omega$, and $\lceil n/d \rceil \geq 2\varpi - 1$. This result noticeably extends the scope of recovery conditions over the domain of practically relevant (quasi-)regular and stochastic carriers. Among the examples are nonstationary sinusoidal and harmonic signals and arbitrary spike-trains with the distance between neighboring spikes at or below $d$ points. Note that $n_s \geq \lceil n/d \rceil$ by the definition of $\mathcal{C}_d$. Hence the relaxation of the strict regularity condition on the $|c_i| = 1$ sample points comes at the expense of a slightly tighter constraint on $n_s$ necessary for exact recovery of $\mathbf{m}$: compare $n_s \geq 2\varpi - 1$ vs. $n_s \geq \varpi + \omega - 1$.[4]

Another important generalization of the recovery conditions comes with the following inequality:

**Proposition II.4.** *Consider* $\mathbf{m} \in \mathcal{M}_\omega$ *and* $\mathbf{c} \in \mathcal{S}_{|\cdot| \leq 1}$. *Take* $n_s \geq 2\varpi - 1$ *sample points of* $\mathbf{s} = \mathbf{m} \circ \mathbf{c}$ *whose indexes are defined as entries of any chosen* $\mathbf{r} \in \mathbb{N}_+^{n_s}$ *with* $r_{i+1} - r_i = n/n_s$ *for every* $i \in \mathcal{I}_{n_s}$. *Then,*

$$\|\mathbf{m} - \hat{\mathbf{m}}\|_2 / \|\mathbf{m}\|_2 \leq \sqrt{1 - \sum_{i=1}^{n_s} s_{r_i}^2 / \sum_{i=1}^{n_s} m_{r_i}^2}. \quad (7)$$

---

[4]This statement is exact and is established as an intermediate result in the proof of *Proposition II.1*.

Hence, if one can find a sequence of at least $2\varpi - 1$ regularly-spaced sample points with $|s_i|$ sufficiently close to $m_i$, then the relative recovery error is close to 0 in terms of (7). This result endows $\hat{\mathbf{m}}$ with the stability to discrepancies from the recovery conditions discussed earlier. At the same time, it provides approximate recovery guarantees for a wider range of stochastic and (quasi-)regular carriers besides those with fairly densely packed $|c_i| = 1$ sample points. Due to the low-pass restriction on $\mathbf{m}$, (7) is expected to hold approximately for an irregular $\mathbf{r} \in \mathbb{N}_+^{n_s}$ with $r_{i+1} - r_i \leq \lceil n/n_s \rceil$ as well.

We finally note that, whereas $\omega$ and $d$ are fixed properties of $\mathbf{m}$ and $\mathbf{c}$, $\varpi$ is a control parameter that must be specified. An appropriate $\varpi$, which satisfies the recovery conditions formulated above, can only be selected by using prior knowledge on $\mathbf{m}$ and $\mathbf{c}$ or found in a supervised learning setup.

### D. Relaxation of the exact minimum-norm requirement

The norm-minimizing property of $\hat{\mathbf{m}}$ in (6) is critical in formulating sharp recovery conditions. However, from a practical point of view, little would be lost if another estimator $\hat{\mathbf{m}}$ with only slightly larger than the minimum norm among all elements of $\mathcal{S}_{\geq|\mathbf{s}|} \cap \mathcal{S}_\varpi$ is used. Thus, we relax (6) to

$$\begin{aligned} \text{find} \quad & \hat{\mathbf{m}} \in \mathcal{S}_{\geq|\mathbf{s}|} \cap \mathcal{S}_\varpi \\ \text{subject to} \quad & \|\hat{\mathbf{m}}\|_2 \simeq \underset{\mathbf{x} \in \mathcal{S}_{\geq|\mathbf{s}|} \cap \mathcal{S}_\varpi}{\arg\min} \|\mathbf{x}\|_2 \end{aligned} \quad (8)$$

To specify the otherwise ambiguous relation operator $\simeq$, we request that $\hat{\mathbf{m}}$ obtained through (8) recovers $\mathbf{m}$ exactly, i.e., is norm-minimizing, for sinusoidal, harmonic, and spike-train carriers covered by *Proposition II.3*. As we see later, this restriction regularizes the numerical algorithms formulated in the present work for sufficiently accurate demodulation well beyond those three classes of $\mathbf{c}$. The advantage brought by the approximation is computational efficiency.

### E. Method of solution

The algorithms that we introduce to solve (6) and (8) in this work fall in the domain of the so-called methods of alternating projections (APs). The defining feature of each AP method is an iterative calculation of a feasible point ($\hat{\mathbf{m}} \in \mathcal{S}_{\geq|\mathbf{s}|} \cap \mathcal{S}_\varpi$ in our case) via alternating metric projections of its current estimate onto the separate constraint sets. Initially proposed by von Neumann for two closed subspaces [27], this approach was later extended to arbitrary closed convex sets of a Hilbert space (see [28] for a review). Various implementations of the AP algorithms exist, featuring different domains of application, rate of convergence, and additional requirements satisfied by the solutions [28], [29].

We provide a rigorous mathematical basis on which the AP algorithms for solving the demodulation problem rely in Suppl. Mat. D, E, F. For a practical comprehension of the material that follows next, it is sufficient to know that:

- The sets $\mathcal{S}_{\geq|\mathbf{s}|}$ and $\mathcal{S}_\varpi$ are closed and convex.
- A metric projection, or simply a projection henceforth, of $\mathbf{z} \in \mathbb{R}^n$ onto a closed convex set $\mathcal{S} \subset \mathbb{R}^n$ is a unique $\mathbf{x}_\mathbf{z} \in \mathcal{S}$ with the smallest distance, i.e., $\|\mathbf{x}_\mathbf{z} - \mathbf{z}\|_2$, from $\mathbf{z}$.

- The projections onto $\mathcal{S}_{\geq|\mathbf{s}|}$ and $\mathcal{S}_{\varpi}$ are respectively achieved by operators

$$\mathbf{P}_{\mathcal{S}_{\geq|\mathbf{s}|}}[\mathbf{z}] = |\mathbf{s}| + (\mathbf{z} - |\mathbf{s}|) \circ \theta(\mathbf{z} - |\mathbf{s}|) \qquad (9)$$

and

$$\mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{z}] = (\mathbf{F^{-1}\,W_{\varpi}\,F})\,\mathbf{z}. \qquad (10)$$

Here, $\theta(\ldots)$ is the Heaviside step function evaluated elementwise. $\mathbf{W}_{\varpi}$ is a diagonal matrix such that $(W_{\varpi})_{ii} = 1$ if $i \in \mathcal{I}_n^{\varpi}$, and $(W_{\varpi})_{ii} = 0$ otherwise.

To emphasize the nature of the underlying numerical algorithms, we name our new approach as AP demodulation.

### F. Relation to other problems and approaches

Demodulation is a counterpart of a widely known and studied problem of blind deconvolution: $\mathbf{s} = \mathbf{m} \circledast \mathbf{c}$. Indeed, both tasks admit the algebraic form of each other in the Fourier domain. Nevertheless, the properties of $\mathbf{m}$ and $\mathbf{c}$ inherent to practical instantiations of amplitude demodulation and blind deconvolution differ significantly. These differences predetermine the need for distinctive strategies to solve the respective tasks, as discussed next.

One of the most powerful convex-programming-based deconvolution approaches, introduced in [30], builds on the assumption that $\mathbf{m}$ and $\mathbf{c}$ belong to known low-dimensional subspaces. There, recovery of $\mathbf{m}$ and $\mathbf{c}$ is achieved by minimizing the nuclear, atomic, $\ell_1$, or $\ell_{2,1}$ norms of their outer product (in the subspace representation) subject to linear measurement constraints of $\mathbf{s}$ [30]–[33]. This scheme successfully solves many practically relevant blind deconvolution cases, such as image deblurring, multipath channel protection, and super-resolution microscopy [30], [33]. However, the low-dimension subspace assumption, a crucial prerequisite of the approach, renders it inapt to deal with realistic carriers in the amplitude demodulation context. Indeed, even a sinusoidal carrier with a fluctuating phase is hardly representable in this frame, not to mention more complex wideband signals met in practice. Moreover, the subspace model of $\mathbf{m}$ and $\mathbf{c}$ does not allow enforcing the amplitude contents to $\mathbf{m}$ exclusively.

Deconvolution problems have also been approached by using AP-like methods [34]–[37]. A general strategy of the existing algorithms is to achieve deconvolution by an iterative refinement of both $\mathbf{m}$ and $\mathbf{c}$ upon the requirement of exact [34], [36] or approximate [35], [37] adherence to the defining equality $\mathbf{s} = \mathbf{m} \circledast \mathbf{c}$ and the support region, intensity range, and spectrum constraints implied on $\mathbf{m}$ and $\mathbf{c}$ or $\mathbf{r} = \mathbf{s} - \mathbf{m} \circledast \mathbf{c}$. These methods differ significantly between themselves. Each of them achieves satisfactory recovery by a judicious combination of specific constraint sets and the iterative scheme adjusted to specific classes of $\mathbf{m}$ and $\mathbf{c}$. The nonconvexity of $\mathcal{C}_d$ and the absence of efficient explicit projections onto this set makes the application of the known deconvolution methods unsuitable to amplitude demodulation. None of the current AP-like deconvolution methods allow assigning the amplitude contents to $\mathbf{m}$ purely either.

We next note that our formulation of the amplitude demodulation problem in Section II-B reveals it as a generalization of the classical task of band-limited signal recovery from true sample points. An AP method known under the name Papoulis-Gerchberg and its variants were successfully applied in the latter setting (see [38] for a review). The differences in the available information on the recoverable signal lead to distinct strategies in algorithmic approaches to these two problems. In particular, the Papoulis-Gerchberg methods rely entirely on known true data. Thus, they are impossible to use for demodulation purposes. The AP algorithms introduced in the present work can be applied in the classical setting. However, not using the available information about the true data makes them inferior to their classical counterparts unless the sample points are fairly uniformly spread, as discussed in Section II-C.

The approach suggested in the present work also has some parallels with the LDC demodulation method by [19]. There, (1) is accompanied by a constraint on the modulator $\mathbf{m}$ expressed as the solution of the quadratic programming problem

$$\begin{aligned} \text{minimize} \quad & \|\mathbf{w} \circ \mathbf{Fm}\|_2^2 + \|\mathbf{m}\|_2^2 \\ \text{subject to} \quad & |s_i| \leq m_i \leq \max[\mathbf{s}] \qquad \forall i \in \mathcal{I}_n, \end{aligned} \qquad (11)$$

where $\mathbf{w}$ denotes the weighting vector. (11) was introduced heuristically, trying to quantify the intuitive notion of the modulator-envelope as a signal wrapping $\mathbf{s}$ from above.

Practical applications suggest the LDC method defined by (1) and (11) being computationally most efficient and precise among all current techniques designed for demodulating signals unreachable to classical algorithms [19], [24]. Thus, we use it as a reference when evaluating the performance of the newly-formulated approach of the present work.

### III. DEMODULATION ALGORITHMS

In this section, we formulate three algorithms representing the core arsenal of the AP approach to demodulation. Simplicity, efficiency, and estimation accuracy of the algorithms are the main aspects under consideration. We refer the reader to Suppl. Mat. F for proofs of all propositions found here.

### A. AP-Basic

We start with the simplest possible AP algorithm, therefore named "AP-Basic" (AP-B).

---

**Algorithm:** AP-Basic (AP-B)

---
1: **Set:** $N_{iter}$, $\epsilon_{tol}$
2: **Initialize:** $i = 0$, $\epsilon^{(0)} = \|\mathbf{s}\|_2/\sqrt{n}$, $\mathbf{m}^{(0)} = |\mathbf{s}|$, $\mathbf{a}^{(0)} = \mathbf{0}$
3: **while** $\epsilon^{(i)} > \epsilon_{tol}$ **and** $i < N_{iter}$ **do**
4: $\quad i = i + 1$
5: $\quad \mathbf{a}^{(i)} = \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{m}^{(i-1)}]$
6: $\quad \mathbf{m}^{(i)} = \mathbf{P}_{\mathcal{S}_{\geq|\mathbf{s}|}}[\mathbf{a}^{(i)}]$
7: $\quad \epsilon^{(i)} = \|\mathbf{m}^{(i)} - \mathbf{a}^{(i)}\|_2/\sqrt{n}$
8: **end while**
9: **Finalize:** $\hat{\mathbf{m}} = \mathbf{m}^{(i)}$

---

Here, $N_{iter}$ stands for the maximum number of algorithm iterations. $\epsilon^{(i)}$ is the infeasibility error at the $i$-th iteration, which is used to control the termination of the algorithm. Specifically, $\epsilon^{(i)}$ measures the distance of the modulator estimate $\mathbf{m}^{(i)} \in \mathcal{S}_{\geq|\mathbf{s}|}$ from $\mathcal{S}_{\varpi}$ and sets a lower bound on
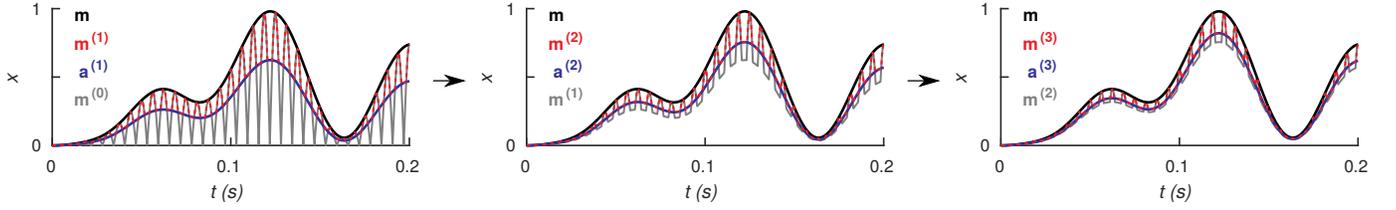
Fig. 1. The first three iterations of the AP-B algorithm applied to an amplitude-modulated sinusoidal signal. $\mathbf{m}$ stands for the real modulator.

the convergence error: $\epsilon^{(i)} \le \|\mathbf{m}^{(i-1)} - \mathbf{m}^{\dagger}\|_2/\sqrt{n}$ (see Suppl. Mat. G). The iterative process is stopped when $\epsilon^{(i)}$ drops to the level of a predetermined threshold $\epsilon_{tol} > 0$ or below. $\epsilon_{tol} \le 0$ would force the completion of all $N_{iter}$ iterations of the algorithm. $\hat{\mathbf{m}}$ denotes the final estimate of the modulator. $\hat{\mathbf{m}}$ arbitrarily close to $\mathcal{S}_{\ge|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$ can be reached if $N_{iter}$ is sufficiently large:

**Proposition III.1.** *A sequence* $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ *formed by the AP-B algorithm for* $\epsilon_{tol} = 0$ *and* $N_{iter} \to +\infty$ *converges to some* $\mathbf{m}^{\dagger} \in \mathcal{S}_{\ge|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$. *The convergence is geometric and monotonic, i.e., there exist* $\gamma > 0$ *and* $0 < r < 1$ *such that* $\|\mathbf{m}^{(i)} - \mathbf{m}^{\dagger}\|_2 \le \gamma \cdot r^i$ *and* $\|\mathbf{m}^{(i+1)} - \mathbf{m}^{\dagger}\|_2 \le \|\mathbf{m}^{(i)} - \mathbf{m}^{\dagger}\|_2$ *for* $i \ge 0$.

It can be shown by example that the AP-B does not always provide minimum-norm estimators $\hat{\mathbf{m}}$. However, it is expected to do so at least approximately if some conditions are met. We clarify this next with the help of Fig. 1, which displays the first three iterations of the AP-B applied to an example signal.

First, note that the starting point $\mathbf{m}^{(0)}$ is elementwise not-higher than the real modulator $\mathbf{m}$ (black). $\mathbf{P}_{\mathcal{S}_{\varpi}}$ maps $\mathbf{m}^{(0)}$ to $\mathbf{a}^{(1)}$, which, by definition of a metric projection, is its best mean-squared-error (MSE) approximation in $\mathcal{S}_{\varpi}$ (blue). By the definition of $\mathcal{S}_{\varpi}$, $\mathbf{a}^{(1)}$ is nearly constant over time windows shorter than $n/(2\pi\varpi)$ points. In general, the best constant MSE estimator of a sample of numbers is its average. Thus, as the best MSE estimator of $\mathbf{m}^{(0)}$ over $\mathcal{S}_{\varpi}$, $\mathbf{a}^{(1)}$ approximates the local average of $\mathbf{m}^{(0)}$ values in a window of $\approx n/(2\pi\varpi)$ points at every moment. If $\varpi \ge \omega$, $\mathbf{c} \in \mathcal{S}_{|\cdot|\le 1}$, and $\approx n/(2\pi\varpi)$ sample points are sufficient to average out the local variations of $\mathbf{c}$, $\mathbf{a}^{(1)}$ is supposed to be proportional to $\mathbf{m}$, at least roughly. The first iteration is completed by the projection of $\mathbf{a}^{(1)}$ back onto $\mathcal{S}_{\ge|\mathbf{s}|}$ to obtain $\mathbf{m}^{(1)}$ (red).

Applying the same reasoning as above, we deduce that, with each iteration, $\mathbf{a}^{(i)}$, and thus $\mathbf{m}^{(i)}$, approaches $\mathbf{m}$ elementwise (see Fig. 1). In general, $\mathbf{m}^{(i)}$ may exceed the level of the real modulator $\mathbf{m}$ over time windows longer than $\ge n/(2\pi\varpi)$ points for higher $i$ before $\mathbf{m}^{\dagger} \in \mathcal{S}_{\ge|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$ is reached. However, as follows from the considerations of the previous paragraph, such segments of $\mathbf{m}^{(i)}$ would be approximately compatible with $\mathcal{S}_{\ge|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$ and would not be considerably affected in subsequent iterations. Hence, $\hat{\mathbf{m}}$ obtained by the AP-B is expected to follow the true sample points of $\mathbf{m}$ tightly. If the number of these points is sufficient, then $\hat{\mathbf{m}} \simeq \mathbf{m}$ as well.

The basis for the above considerations is laid by the fact that they are exact for some important types of carriers:

**Proposition III.2.** *Consider* $\mathbf{m} \in \mathcal{M}_{\omega}$ *and* $\mathbf{c} \in \mathcal{C}_d$ *with* $|c_j| = \sum_{k=1}^{n/\nu}(\tilde{c}_{\nu\cdot k} \cdot e^{\imath 2\pi\nu(k-1)(j-1)/n})$, *where* $\tilde{c}_{\nu\cdot k} \in \mathbb{C}$ *and* $n/\nu \in \mathbb{N}$. *If* $\varpi \ge \omega$ *and* $\nu \ge \varpi + \omega - 1$, *then a sequence* $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ *formed by the AP-B algorithm for* $\epsilon_{tol} = 0$ *and* $N_{iter} \to +\infty$ *converges to* $\mathbf{m}$.

Among others, *Proposition III.2* encompasses the sinusoidal, harmonic, and regular spike-train carriers covered by *Proposition II.3*. Thus, in these cases, AP-B satisfies the minimum-norm property, i.e., provides $\hat{\mathbf{m}}$ that converges to a solution of (8). The condition $\nu \ge \varpi + \omega - 1$ in *Proposition III.2* plays the role of the inequality $n/d \ge \varpi + \omega - 1$ in *Proposition II.3*.

### B. AP-Accelerated

One of the potential weak points of AP algorithms based on pure alternating projections onto convex sets, like the AP-B, is relatively slow convergence [39]–[41]. Indeed, despite the geometric nature of the convergence, the actual number of iterations necessary to reach a specific error level may be arbitrarily large if the factor $r$ in $\|\mathbf{m}^{(i)} - \mathbf{m}^{\dagger}\|_2 \le \gamma \cdot r^i$ is sufficiently close to 1. To address this issue, various accelerated AP schemes have been suggested for specific classes of the constraint sets [39], [42], [43]. Here, we propose a parameter-free accelerated version of the AP-B algorithm specifically designed for the demodulation problem. We refer to it as "AP-Accelerated" (AP-A).

---

**Algorithm: AP-Accelerated (AP-A)**

1: **Set:** $N_{iter}$, $\epsilon_{tol}$
2: **Initialize:** $i = 0$, $\epsilon^{(0)} = \|\mathbf{s}\|_2/\sqrt{n}$, $\mathbf{m}^{(0)} = |\mathbf{s}|$, $\mathbf{a}^{(0)} = \mathbf{0}$
3: **while** $\epsilon^{(i)} > \epsilon_{tol}$ **and** $i < N_{iter}$ **do**
4:     $i = i + 1$
5:     $\mathbf{b}^{(i)} = \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{m}^{(i-1)} - \mathbf{a}^{(i-1)}]$
6:     $\lambda = \|\mathbf{m}^{(i-1)} - \mathbf{a}^{(i-1)}\|_2^2/\|\mathbf{b}^{(i)}\|_2^2$
7:     $\mathbf{a}^{(i)} = \mathbf{a}^{(i-1)} + \lambda \cdot \mathbf{b}^{(i)}$
8:     $\mathbf{m}^{(i)} = \mathbf{P}_{\mathcal{S}_{\ge|\mathbf{s}|}}[\mathbf{a}^{(i)}]$
9:     $\epsilon^{(i)} = \|\mathbf{m}^{(i)} - \mathbf{a}^{(i)}\|_2/\sqrt{n}$
10: **end while**
11: **Finalize:** $\hat{\mathbf{m}} = \mathbf{m}^{(i)}$

---

**Proposition III.3.** *A sequence* $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ *formed by the AP-A algorithm for* $\epsilon_{tol} = 0$ *and* $N_{iter} \to +\infty$ *converges to some* $\mathbf{m}^{\dagger} \in \mathcal{S}_{\ge|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$. *The convergence is monotonic, i.e.,* $\|\mathbf{m}^{(i+1)} - \mathbf{m}^{\dagger}\|_2 \le \|\mathbf{m}^{(i)} - \mathbf{m}^{\dagger}\|_2$ *for* $i \ge 0$.

Note that $\lambda > 1$ except when $\mathbf{P}_{\mathcal{S}_{\varpi}}$ is the identity operator, i.e., the trivial case of $\mathbf{m} = |\mathbf{s}|$. Indeed, it follows from the definition of $\mathbf{P}_{\mathcal{S}_{\varpi}}$ [see (10)] and the unitary property of $\mathbf{F}$ that

$\|\mathbf{m}^{(i-1)} - \mathbf{a}^{(i-1)}\|_2^2 > \|\mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i-1)} - \mathbf{a}^{(i-1)}]\|_2^2 = \|\mathbf{b}^{(i)}\|_2^2$, if $\mathbf{P}_{\mathcal{S}_\varpi}$ is not the identity operator. It is easy to see that $(\mathbf{a}^{(i)} - \mathbf{a}^{(i-1)}) = \lambda \cdot (\mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i-1)}] - \mathbf{a}^{(i-1)})$ in the above algorithm. Moreover, if $\lambda$ is fixed to 1 by force, the AP-A and AP-B algorithms become identical. Therefore, the AP-A produces increments from $\mathbf{a}^{(i-1)}$ to $\mathbf{a}^{(i)}$ that are scaled up compared with those that were obtained by applying the AP-B algorithm for the same iterations.

To understand the working principle of the AP-A better, recall that $\mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}]$, and thus $\mathbf{a}^{(i)}$, are nearly constant over time windows consisting of $< n/(2\pi\varpi)$ points (see Section III-A). For a semiquantitative analysis, we can assume that this holds exactly. Let us denote a segment of $(\mathbf{m}^{(i-1)} - \mathbf{a}^{(i-1)})$ restricted to such a window by $\mathbf{z}$. Then, $\mathbf{b}^{(i)}$ defined in the same window is just $(l^{-1} \cdot \sum_{j=1}^l z_j) \cdot \mathbf{1}$, and $\|\mathbf{m}^{(i-1)} - \mathbf{a}^{(i-1)}\|_2^2$ corresponds to $\sum_{j=1}^l z_j^2$, where, $l = \lfloor n/(2\pi\varpi) \rfloor$. Consequently, $\lambda \cdot \mathbf{b}^{(i)}$, i.e., the increment from $\mathbf{a}^{(i-1)}$ to $\mathbf{a}^{(i)}$, is given by $\left(\sum_{j=1}^l z_j^2 / \sum_{j=1}^l z_j\right) \cdot \mathbf{1}$. It follows from $\mathbf{m}^{(i-1)} = \mathbf{P}_{\mathcal{S}_{\geq|\mathbf{s}|}}[\mathbf{a}^{(i-1)}]$ that $\mathbf{z}$ is elementwise nonnegative. Therefore, $\left(\sum_{j=1}^l z_j^2 / \sum_{j=1}^l z_j\right) \leq \max[\mathbf{z}]$. However, $\max[\mathbf{z}]$ corresponds to the difference between the real modulator and $\mathbf{a}^{(i-1)}$ in the considered time window, at least approximately, if $\lceil n/d \rceil \geq 2\varpi - 1$. Thus, while up-scaling $\mathbf{a}^{(i)} - \mathbf{a}^{(i-1)}$ at each iteration to accelerate the convergence, the AP-A also ensures that $\mathbf{a}^{(i)}$ stays approximately within the bounds of the real modulator $\mathbf{m}$. This property ensures that $\hat{\mathbf{m}}$ tightly follows $\mathbf{m}$ if the same conditions as required by the AP-B are met.

We further note that $\left(\sum_{j=1}^l z_j^2 / \sum_{j=1}^l z_j\right) = \max[\mathbf{z}]$, i.e., $\mathbf{a}^{(i)}$ reaches $\mathbf{m}$ in a single iteration, if all but one element of $\mathbf{z}$ are equal to zero. Importantly, approximately this situation is faced in reality, as illustrated in Fig. 1. Specifically, with increased $i$, $(\mathbf{m}^{(i-1)} - \mathbf{a}^{(i-1)})$ becomes mainly flat with only a few separate elements considerably above 0 over time windows shorter than $n/(2\pi\varpi)$ points. For comparison, the analogous increment from $\mathbf{a}^{(i-1)}$ to $\mathbf{a}^{(i)}$ is moderate and equals only $\max[\mathbf{z}]/l$ in the case of the AP-B method. These considerations explain the substantial speed-up provided by the AP-A algorithm in practice. They also reveal that any additional acceleration steps in the AP-A would result in overscaled $\hat{\mathbf{m}}$, hence reducing the demodulation accuracy.

The AP-A algorithm repeats the AP-B in terms of exact recovery guarantees of *Proposition III.2*:

**Proposition III.4.** *Consider* $\mathbf{m} \in \mathcal{M}_\omega$ *and* $\mathbf{c} \in \mathcal{C}_d$ *with* $|c_j| = \sum_{k=1}^{n/\nu} (\tilde{c}_{\nu \cdot k} \cdot e^{\imath 2\pi\nu(k-1)(j-1)/n})$, *where* $\tilde{c}_{\nu \cdot k} \in \mathbb{C}$ *and* $n/\nu \in \mathbb{N}$. *If* $\varpi \geq \omega$ *and* $\nu \geq \varpi + \omega - 1$, *then a sequence* $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ *formed by the AP-A algorithm for* $\epsilon_{tol} = 0$ *and* $N_{iter} \to +\infty$ *converges to* $\mathbf{m}$.

This result substantiates the semiquantitative argumentation of the AP-A convergence properties provided above and establishes the respective $\hat{\mathbf{m}}$ as a numerical solution of (8).

### C. AP-Projected

As argued above, the AP-A and AP-B algorithms produce modulator estimates that are expected to tightly follow the original $\mathbf{m}$ if the conditions analogous to those discussed in Section II-C are met. These estimates, however, do not always hold the minimum-norm property (6). A classical AP scheme that guarantees minimum-norm solutions is known under the name of Dykstra [44], [45]. In particular, Dykstra's algorithm calculates the projection of a point in $\mathbb{R}^n$ onto the feasible set. Thus, by choosing an appropriate initial condition, the solution with a minimized norm can be obtained (see *Proposition III.5* next and its proof in Suppl. Mat. F). We consider a version of this algorithm adapted to solve the demodulation problem and call it "AP-Projected" (AP-P).

---

**Algorithm:** AP-Projected (AP-P)

1: **Set:** $N_{iter}$, $\epsilon_{tol}$
2: **Initialize:** $i = 0$, $\epsilon^{(0)} = \|\mathbf{s}\|_2/\sqrt{n}$, $\mathbf{m}^{(0)} = \mathbf{c}^{(0)} = |\mathbf{s}|$
3: **while** $\epsilon^{(i)} > \epsilon_{tol}$ **and** $i < N_{iter}$ **do**
4:     $i = i + 1$
5:     $\mathbf{a}^{(i)} = \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i-1)}]$
6:     $\mathbf{m}^{(i)} = \mathbf{P}_{\mathcal{S}_{\geq|\mathbf{s}|}}[\mathbf{a}^{(i)} - \mathbf{c}^{(i-1)}]$
7:     $\mathbf{c}^{(i)} = \mathbf{m}^{(i)} - (\mathbf{a}^{(i)} - \mathbf{c}^{(i-1)})$
8:     $\epsilon^{(i)} = \sqrt{(\|\mathbf{m}^{(i-1)} - \mathbf{a}^{(i)}\|_2^2 + \|\mathbf{m}^{(i)} - \mathbf{a}^{(i)}\|_2^2)/(2 \cdot n)}$
9: **end while**
10: **Finalize:** $\hat{\mathbf{m}} = \mathbf{m}^{(i)}$

---

**Proposition III.5.** *A sequence* $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ *formed by the AP-P algorithm for* $\epsilon_{tol} = 0$ *and* $N_{iter} \to +\infty$ *converges to a unique* $\mathbf{m}^\dagger \in \mathcal{S}_{\geq|\mathbf{s}|} \cap \mathcal{S}_\varpi$ *such that* $\|\mathbf{m}^\dagger\|_2 \leq \|\mathbf{x}\|_2$ *for every* $\mathbf{x} \in \mathcal{S}_{\geq|\mathbf{s}|} \cap \mathcal{S}_\varpi$. *The convergence is monotonic, i.e.,* $\|\mathbf{m}^{(i+1)} - \mathbf{m}^\dagger\|_2 \leq \|\mathbf{m}^{(i)} - \mathbf{m}^\dagger\|_2$ *for* $i \geq 0$.

The AP-P differs from the AP-B in that, before projecting a point onto $\mathcal{S}_{\geq|\mathbf{s}|}$, an increment produced by the projection onto this set in the previous iteration is subtracted. This correction may cause the infeasibility error $\|\mathbf{m}^{(i)} - \mathbf{a}^{(i)}\|_2/\sqrt{n}$ estimated after projecting onto $\mathcal{S}_{\geq|\mathbf{s}|}$ to drop to zero intermittently before the final solution is reached, making it an inappropriate option as the stopping criterion. Hence, in contrast to the AP-B and AP-A algorithms, we defined the $\epsilon$ for the AP-P as a combination of the infeasibility errors evaluated after projecting onto both sets $\mathcal{S}_\varpi$ and $\mathcal{S}_{\geq|\mathbf{s}|}$ at each iteration (see line 8 above). This error measure is strictly positive and converges to zero when $N_{iter} \to +\infty$ [46].

The understanding of the convergence rate of Dykstra's scheme is limited. It was shown that the convergence is geometric for an intersection of half-spaces [47], [48]. Nevertheless, no equivalent result exists for other convex sets. Moreover, it was demonstrated that the convergence rate of this algorithm may depend on the initial conditions and may be considerably slower than that of AP algorithms based on pure projections [49].

### D. Computational complexity

Except for the projection operator $\mathbf{P}_{\mathcal{S}_\varpi}$, each iteration of the three formulated AP algorithms relies on vector addition, scalar product, and value update. These are linear in the number of sample points. The operator $\mathbf{P}_{\mathcal{S}_\varpi}$ can be easily implemented by using the direct and inverse fast Fourier transforms (FFTs) and setting the relevant elements of the signal to
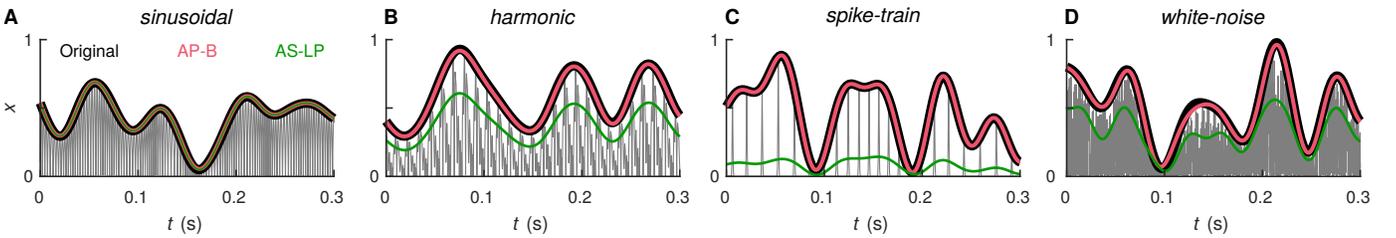
Fig. 2. Typical examples of the test signals (gray), featuring nonstationary sinusoidal (**A**), harmonic (**B**), spike-train (**C**), and stationary white-noise (**D**) carriers, and their modulators obtained by using the AP-B (red) and AS-LP (green) algorithms. The signals are represented by their absolute values here. The predefined modulators are shown in black.

zero in the Fourier domain. The current state-of-the-art FFT algorithms have an $\mathcal{O}(n \log n)$ time complexity [50], which, thus, sets the overall time complexity of the AP algorithms introduced in this work. Our numerical experiments suggest that the convergence speed in terms of iteration number is independent of the signal length (see Suppl. Mat. M).

## IV. PERFORMANCE TESTS

To evaluate the AP algorithms introduced above, we compared their performance with the AS and LDC demodulation approaches when applied to infer the modulator of predefined synthetic test signals. The LDC approach was implemented by using two state-of-the-art quadratic programming solvers: Gurobi (v8.1.1) [51] and OSQP (v0.6.0) [52]. The AS demodulation was achieved by using the FFT-based approach [53]. In that case, we additionally low-pass filtered the obtained modulator estimate with $\mathbf{P}_{\mathcal{S}_\varpi}$ to regularize it. We refer to this modified demodulation scheme as AS-LP.

### A. Test signals

The test signals were composed as products of a modulator and a carrier: $\mathbf{s} = \mathbf{m} \circ \mathbf{c}$. Four types of $\mathbf{c}$, approximating basic building blocks of real-world signals, were used: nonstationary *sinusoidal*, *harmonic*, and *spike-train*, as well as stationary *white-noise* (see, respectively, (176), (180), (184), and (188) in Suppl. Mat. I). The former two were combined with modulators of nonstationary *Gaussian* origin, while the latter two types of carriers were paired with the so-called *maximally-uniformly distributed* modulators (see, respectively, (160)−(162) and (163)−(168) in Suppl. Mat. I).

In all cases, modulator and carrier pairs were selected to meet the core recoverability condition $\lceil n/d \rceil \geq 2\varpi - 1$, at least approximately. The remaining parameters of $\mathbf{m}$ and $\mathbf{c}$ (see Suppl. Mat. I) were chosen to imitate realistic conditions as much as possible. For example, in all cases, signals were taken as segments of longer time series, and thus, were not $n$-periodic. The center frequencies of the sinusoidal and harmonic carriers were set so that only sample points with $|c_i| \approx 1$ rather than $|c_i| = 1$ were available.

### B. Performance evaluation

Demodulation performance was evaluated by using two complementary measures: 1) error of the modulator estimate, $E_m = \|\mathbf{m} - \hat{\mathbf{m}}\|_2 / \|\mathbf{m}\|_2$; and 2) execution time of the algorithm on the computer, $T_{\text{cpu}}$. We evaluated the AP and LDC

algorithms in the mode when $T_{\text{cpu}}$ depends on the total number of sample points but not on the effective degrees of freedom. This choice made the results general, independent of a selected cutoff frequency $\varpi$. To insure against outliers, we averaged $E_m$ and $T_{\text{cpu}}$ over ten independent signal realizations.

Execution of the AP and LDC algorithms is controlled by a set of metaparameters whose choice influences the output. Therefore, we aimed for the Pareto fronts, not separate points, in the $(E_m, T_{\text{cpu}})$ plane. Due to the computing speed limitations inherent to the LDC approach, we had to exploit signal decomposition into separate fragments for this analysis. In particular, signals were split into segments that were demodulated separately and then put together [19]. This allowed achieving a linear growth in the computation time with the total length of the signal, and hence, speeding up the calculations. After identifying the optimal control-parameter combinations, we compared all methods by demodulating whole signals.

Sets of the demodulation control parameters that we considered for the Pareto optimality analysis, including those defining the signal splitting, are provided in Suppl. Mat. J. Details on the execution of the performance tests on a computer can be found in Suppl. Mat. K.

### C. Results

Fig. 2 shows representative fragments of the test signals from all four classes (gray) and their modulator estimates obtained by using the AS-LP (green) and AP-B (red) algorithms. Whereas the AP-B allows obtaining high-quality estimates $\hat{\mathbf{m}}$ in all four cases, the AS-LP does so only for sinusoidal signals.

Results of the performance evaluation in the form of Pareto fronts in the $(E_m, T_{\text{cpu}})$ plane for $n = 2^{15}$ are displayed in Fig. 3 A–D. Panels E–H of the same figure show $T_{\text{cpu}}$ vs. $n$ relations derived by using no window splitting. A closer analysis of these data reveals the following:

1) The AP algorithms feature lower bounds on the demodulation error $E_m$ than the LDC method (Fig. 3 A–D).
2) The AP algorithms are up to five orders of magnitude faster than their LDC counterparts for achieving the same $E_m$ when optimal signal window splitting is used (Fig. 3 A–D). The difference is even more pronounced when no window splitting is assumed (Fig. 3 E–H). For example, to process a 1 s length signal sampled at 16 kHz, the LDC needs $10^4$ s of CPU time, in contrast to $10^{-3}$ s taken by the AP-A.
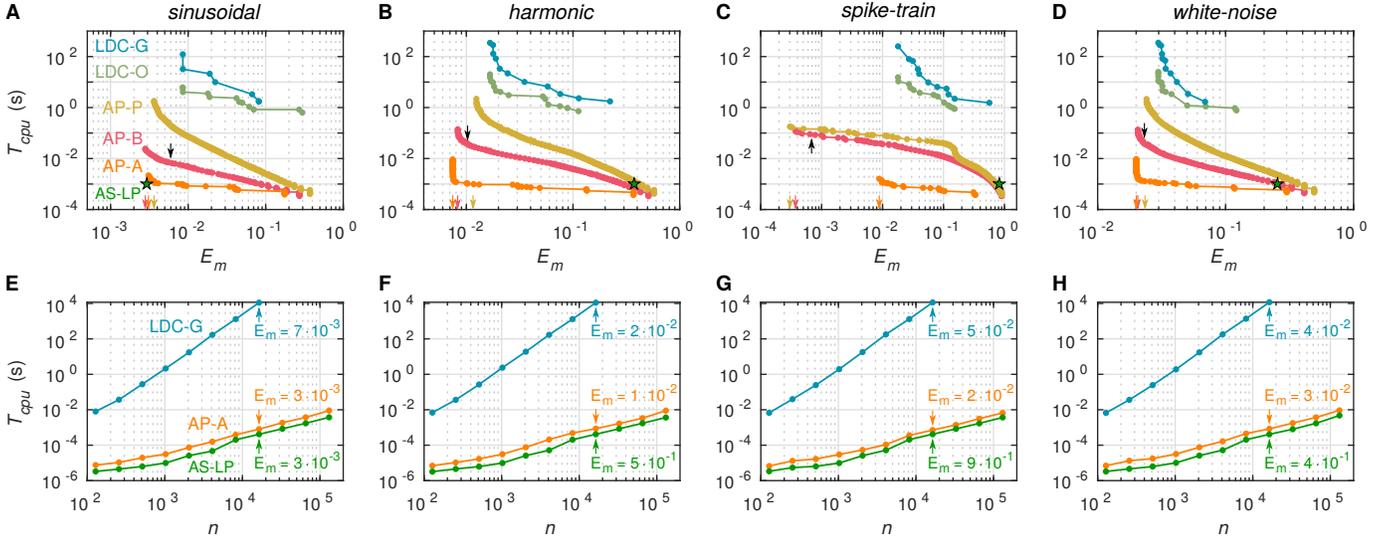
Fig. 3.  Performance evaluation. **A–D**: Pareto fronts in the $(E_m, T_{\rm cpu})$ plane for different demodulation algorithms applied to the four different types of test signals when window splitting is used. Green stars mark the results of the AS-LP method. Color arrowheads point to the lower bounds on $E_m$ for the respective AP algorithms. Black arrowheads indicate demodulation error $E_m$ values of the AP-B algorithm calculated locally for signal windows shown in Fig. 2. **E–H**: Dependence of the demodulation time $T_{\rm cpu}$ on the signal length $n$ at $\epsilon_{tol} = 10^{-3}$ when window splitting is not exploited. $E_m$ values in the legends correspond to demodulation results at $n = 2^{14} \approx 1.6 \cdot 10^4$.

3) $T_{\rm cpu}$ varies substantially (up to three orders of magnitude) even between different AP algorithms (Fig. 3 A–D). The AP-A ranks as the fastest, and the AP-P as the slowest one for all tested signals.

4) Despite the differences in $T_{\rm cpu}$, all AP algorithms feature similar lower bounds on $E_m$, except the spike-train signals, when the AP-B and AP-P can noticeably surpass the AP-A on the relative scale (Fig. 3 A–D). Nevertheless, on the absolute scale, the AP-A still performs reasonably well.

5) For all tested signals, the AP-A algorithm outperforms the AS-LP-based demodulation in the sense that it can achieve the same or smaller errors with the same $T_{\rm cpu}$ (Fig. 3 A–D). Moreover, compared with the AS-LP, AP algorithms exhibit much lower bounds on $E_m$.

6) Even without the window splitting (when the highest demodulation accuracy is attained), the AP-A algorithm takes only 2–3 times longer than the AS-LP method (Fig. 3 E–H).

We found that the decrease in $E_m$ along the Pareto fronts is mainly determined by the increase in the demodulation window size. In particular, the lower bounds on $E_m$ are achieved by the particular algorithms when the signal is demodulated using no window splitting. The relatively lower precision of the AP-A algorithm compared with AP-B and AP-P in the case of nonstationary spike-trains can be reduced to its acceleration mechanism. Indeed, in the AP-A, upscaling of iterates $\mathbf{a}^{(i)}$ is effectively based on the averaging of $(\mathbf{m}^{(i)} - \mathbf{a}^{(i)})$ over a window of length $\approx n/(2\pi\varpi)$ at each sample point. The precision of these estimates is more vulnerable to deviations from the exact recovery conditions for sparse carriers.

As can be expected, the high accuracy of modulator estimates achieved by the AP algorithms implies the high quality of carrier predictions $\hat{\mathbf{c}} = \mathbf{s}/\hat{\mathbf{m}}$ (see Suppl. Mat. L and Fig. 13

therein). The AP approach leaves the AS-LP behind in terms of carrier estimation for all four signal types considered (see Suppl. Mat. L). When applicable, the inferred $\hat{\mathbf{c}}$ can be further frequency-demodulated by using dedicated techniques (see [1], [54], and references given there).

The impressive performance of the AP-A algorithm in terms of $E_m$, $E_c$, and $T_{\rm cpu}$ makes it an ideal candidate for amplitude demodulation of a wide range of signals. Its AP-B and AP-P counterparts can be used instead if higher precision is needed in specific cases, as illustrated by the spike-train signals above.

## V. Convergence Tests

To clarify the differences between the $T_{\rm cpu}$ estimates of the three AP algorithms and understand the relationship between the demodulation and infeasibility errors, we performed a convergence analysis with the test signals from the previous section. The simulation results for fixed $n = 2^{15}$ using no window splitting are summarized in Fig. 4. A closer inspection uncovers the following:

1) The convergence rates in terms of both $\epsilon$ and $E_m$ parallel the differences in the computing speed of different AP algorithms. Among them, the fastest is the AP-A, which reaches any given $\epsilon$ or $E_m$ level with the smallest number of iterations. The AP-P algorithm is the slowest one.

2) The AP-A algorithm converges in a finite number of iterations $(< 30)$ for all types of test signals studied. In particular, it requires only $\leq 5$ iterations to reach the plateau level of the demodulation error $E_m$. This fact explains the extraordinary computational efficiency of the AP-A documented in Section IV.

3) Differently from the convergence error $\|\mathbf{m}^{(i)} - \mathbf{m}^\dagger\|_2/\sqrt{n}$, the dependence of $E_m^{(i)}$ on $i$ can be non-monotonic if $\mathbf{m}^\dagger$ is not strictly equal to $\mathbf{m}$ (Fig. 4 E, G). Then, $E_m^{(i)}$ starts growing with increased $i$ after reaching
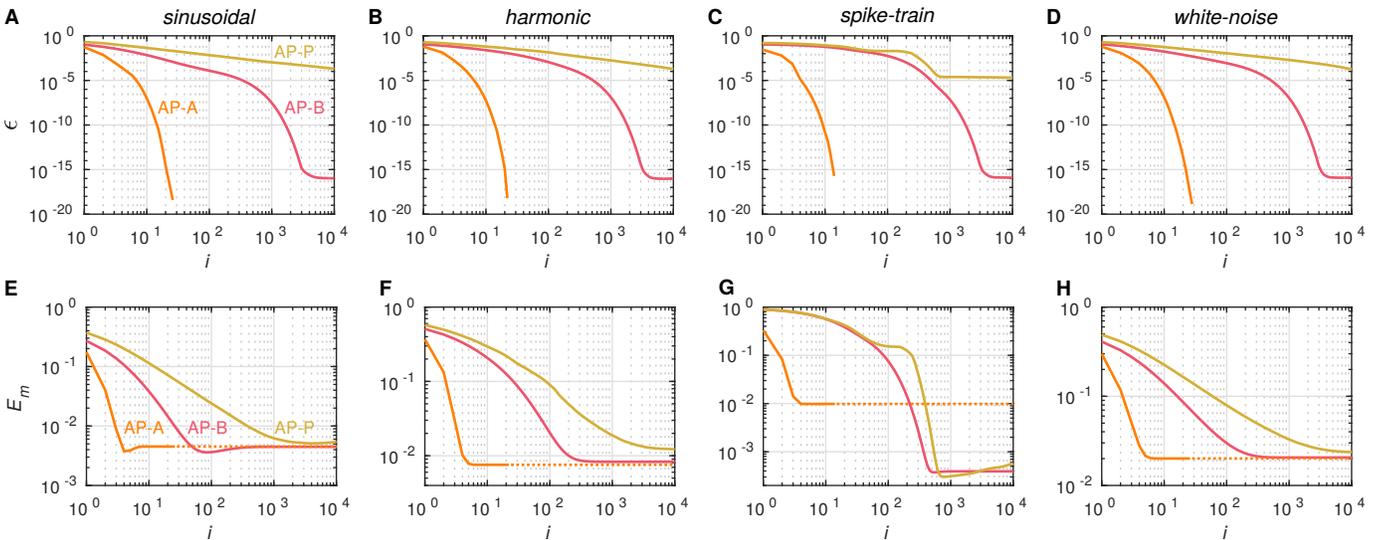
Fig. 4. Convergence analysis of the AP algorithms. **A–D**: Dependence of the infeasibility error $\epsilon$ on the iteration number $i$ for the AP-B, AP-A, and AP-P algorithms applied to the four different types of test signals with $n = 2^{15}$ and no window splitting. **E–H**: Analogous plots to A–D made for the demodulation error $E_m$ instead of $\epsilon$. Dotted lines show hypothetical $E_m$ values that would be obtained if we continued the AP-A iterations after reaching the final solution.

the minimum point. However, this growth is mild and of no practical importance as long as $\mathbf{m}^\dagger \approx \mathbf{m}$, i.e., $\hat{\mathbf{m}} \approx \mathbf{m}$.

The results shown in Fig. 4 represent only signals of fixed length ($n = 2^{15}$ sample points). Additional simulations suggested no dependence on $n$ (see Suppl. Mat. M).

## VI. ROBUSTNESS TESTS

The pivotal condition for successfully separating the modulator-carrier information of a given signal by our approach is $\lceil n/d \rceil \geq 2\omega - 1$. In practice, this requirement is not necessarily met. Hence, the choice of a particular demodulation method must be guided not only by the algorithmic efficiency but also robustness to deviations from the ideal recovery conditions. To shed light on this aspect, we considered demodulation of the test signals from Section IV-A corrupted by a multiplicative Bernoulli-$\{0,1\}$ noise. In this setup, sample points, including the decisive $|s_i| = m_i$, are eliminated with the probability of "0" elements in the noise ($P(0)$), effectively decreasing the value of $\lceil n/d \rceil$.

We found that all three AP algorithms considered in this work show a similar degree of robustness to increased $P(0)$ (see Fig. 5). Only in the case of sinusoidal signals, the AP-A is slightly inferior to the AP-B and AP-P. Interestingly, the advantage of the AP-B and AP-P over the AP-A in the case of spike-train signals discussed in Section IV-C disappears in the presence of even small distortions (see Fig. 5 C). The differences in the $E_m$ vs. $P(0)$ relations seen in Fig. 5 A–D are predetermined by different densities of $|c_i| \simeq 1$ points inherent to each carrier type. Analogous results to those shown in Fig. 5 A–D are obtained when considering carrier recovery via $\hat{\mathbf{c}} = \mathbf{s}/\hat{\mathbf{m}}$ (see Fig. 14 in Suppl. Mat. L).

In contrast to the AP approach, the AS-based demodulation is highly vulnerable to missing sample points, and hence, to decreased $\lceil n/d \rceil$ (Fig. 5 A–D). Even for sinusoidal signals,

which the AS and AS-LP are specially designed for, the zeroing of data points leads to a rapid decline in demodulation quality (Fig. 5 A, E).

The robustness to missing sample points endows the AP demodulation method with a highly valuable practical advantage. In particular, it can be exploited in real-world situations when: 1) the sampling rate is low; 2) some segments of the signal values are lost; 3) some sample points are corrupted by noise such that the level of these points can be reduced below the real modulator by low-pass filtering or explicitly identifying them. In this context, the PAD and LDC demodulations compare to the AP approach by construction [23].

## VII. HIGH-LEVEL PROPERTIES

As emphasized in Section II, different demodulation methods can be derived by requesting adherence of the inferred modulators and carriers to a set of particular properties. Typically, various combinations that consist of a few out of many reasonable requirements are sufficient for unique demodulation formulations. However, some of these requirements are inconsistent with each other, making virtually all classical demodulation approaches fail to satisfy one or another essential condition [9], [23], [25]. For example, the AS demodulation method may return an unbounded modulator estimate for a bounded signal [25].

The AP approach formulated in this work is compatible with the following high-level requirements, which have crystallized as inseparable from the notion of proper amplitude demodulation with time [19], [23, Section 3.5.2]:

- **Boundedness:** The modulator and carrier of a bounded signal are bounded. In particular, it is required that $-\infty < \hat{m}_i < +\infty$ and $-1 \leq \hat{c}_i \leq 1$ for every $i \in \mathcal{I}_n$. In the case of the AP approach, the boundedness of the modulator is guaranteed by the convergence of the AP algorithms. The boundedness of the carrier then follows from the constraint $\hat{m}_i \geq |s_i|$ and the fact that $\hat{\mathbf{c}} = \mathbf{s} \circ \hat{\mathbf{m}}^{-1}$.
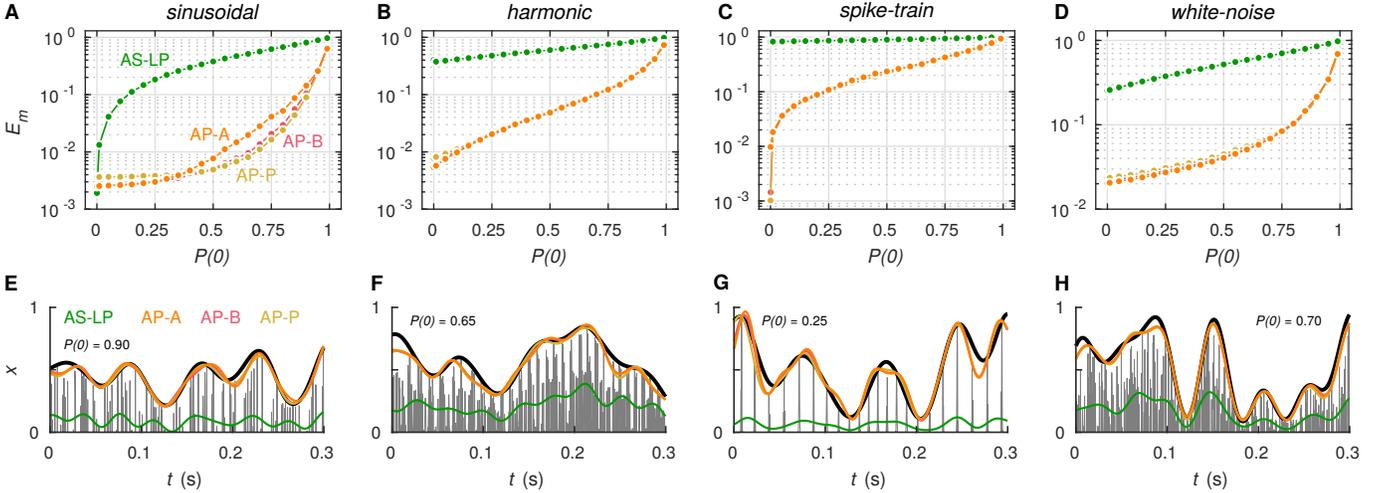
Fig. 5. Robustness evaluation. **A**–**D**: Dependence of the demodulation error $E_m$ on $P(0)$ (the probability of missing points) for the four types of test signals and different AP algorithms at $\epsilon_{tol} = 10^{-4}$ (color coding). **E**–**H**: Representative examples of demodulation at various $P(0)$ levels for the test signals from A–D. Color code: gray – the absolute-value signal, black – the original modulator, color – modulators inferred by different algorithms.

- **Scale covariance:** The modulator and carrier of a scaled signal are equal to the modulator and carrier obtained from the original signal and then scaled by the same amount. The adherence of the AP approach to this condition follows from two facts. First, projection operators $\mathbf{P}_{\mathcal{S}_{\geq|\mathbf{s}|}}$ and $\mathbf{P}_{\mathcal{S}_{\varpi}}$ are homogeneous with degree 1, i.e., $\mathbf{P}_{\mathcal{S}}[\alpha \cdot \mathbf{s}] = \alpha \cdot \mathbf{P}_{\mathcal{S}}[\mathbf{s}]$. Second, each iteration of the AP algorithms can be expressed as a weighted sum of these projections with the weights independent of the scale.
- **Smoothness:** The modulator of a bounded signal in its continuous-time representation is smooth. Because we use a discrete-time representation, this requirement has to be adjusted. In particular, let us denote by $\hat{m}'_t$ and $\hat{m}'_{t+\Delta t}$ the finite-difference approximations of the modulator's time-derivatives of any order at two subsequent time points: $t$ and $t+\Delta t$. Then, we require that, for any $\epsilon > 0$, there exists a $\delta > 0$ such that $|\hat{m}'_{t+\Delta t} - \hat{m}'_t| < \epsilon$ when $|\Delta t| < \delta$. The AP approach satisfies this requirement through the boundedness of the modulator and the bandwidth constraint set by $\mathcal{S}_{\varpi}$ on it.
- **Idempotence:** Information associated with the qualities of modulators and carriers is fully separated. Specifically, demodulation reapplied to an estimated modulator (carrier) must return the same modulator (carrier). The AP approach satisfies the idempotence requirement for the modulator exactly. Indeed, when any AP algorithm is applied to its final solution $\hat{\mathbf{m}} = \mathbf{m}^{\dagger}$, the latter is recognized as the final solution again after the first new iteration by construction. Regarding the carrier, the idempotence holds whenever the recovery conditions discussed in Section II-C are met. That is because, in those cases, $\hat{\mathbf{c}}$ resulting from the first demodulation contains a sufficient number of $|\hat{c}_i| = 1$ points to uniquely define the $\hat{\mathbf{m}} = \mathbf{1}$ as the norm-minimizing element of $\mathcal{S}_{\geq|\hat{\mathbf{c}}|} \cap \mathcal{S}_{\varpi}$. If the recovery conditions are met only approximately, we expect no marked deviations from the idempotence condition (see Fig. 16 in Suppl. Mat. N).

By fulfilling the above requirements, the AP approach parallels the methods of PAD and LDC demodulation [19], [23]. In this sense, all of them outperform the classical techniques.

## VIII. DEMODULATION OF SPEECH SIGNALS

Amplitude demodulation is of central importance in various tasks of processing and analysis of speech signals. Application-wise, this procedure is used in hearing restoration [10], [55], speech recognition [2], [56], [57], and source separation [58], [59]. On the theory side, amplitude demodulation is exploited in neurophysiological and psychophysical studies of auditory information processing in the brain [11], [12], [60], [61]. Depending on the problem, demodulation of either narrow subband [56], [58], intermediate subband [10], [11], or whole wideband signal [12], [62] is needed. In all these cases, modulators and carriers convey the information about specific aspects of speech, e.g., semantic meaning, associated emotion, or speaker identity, that need to be extracted.

In this section, we apply the newly-introduced AP approach to speech demodulation to further demonstrate its potential. To represent the range of possible real-world situations, we consider two limiting signal types: 1) a narrow subband component of a signal obtained by a standard auditory ERB filter [63]; and 2) the original wideband signal.

### A. Direct demodulation

By construction, the output of auditory ERB filters occupies a frequency subband whose width $\Delta$ is much smaller than its center frequency $f_c$ [63]. The resulting signal is an amplitude- and phase-modulated sinusoidal $\mathbf{s} = \mathbf{m} \circ \sin(2\pi f_c \mathbf{t} + \boldsymbol{\varphi})$, with most of the energies of $\mathbf{m}$ and $\boldsymbol{\varphi}$ residing in the frequency interval $[0, \Delta]$ [64]. Hence, by the recovery conditions of the AP approach (see Section II-C), setting the cutoff frequency $\varpi$ between $\Delta$ and $f_c$ necessarily results in accurate estimates $\hat{\mathbf{m}}$ and $\hat{\mathbf{c}}$. In particular, note that the local maximums of $|\mathbf{s}|$ correspond to the true sample points $|s_i| = m_i$. Thus, the high quality of demodulation is visually conveyed by a tight match
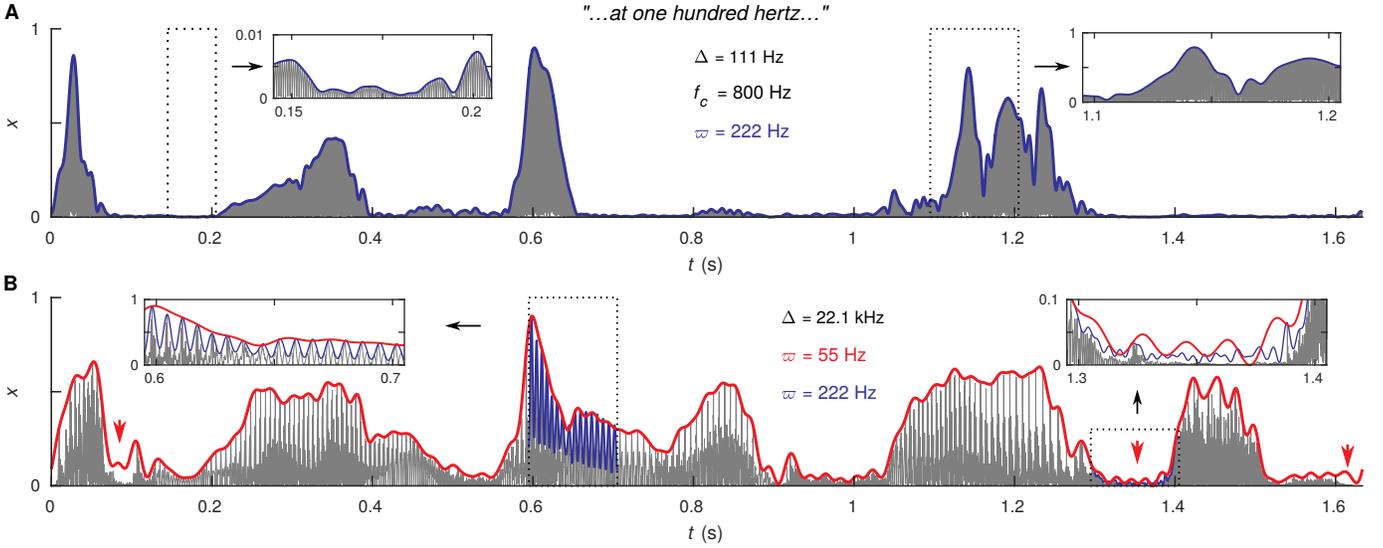
Fig. 6. Direct demodulation of speech signals. **A**: A band-pass-filtered signal of an utterance *"...at one hundred hertz ..."* by a female speaker; the signal was obtained with an equivalent rectangular bandwidth filter of the cochlea centered at $f_c = 800$ Hz and $\Delta = 111$ Hz. **B**: The original signal of the utterance used in panel A (full bandwidth of 22.1 kHz). In both panels, gray color marks the absolute-value version of the signals considered for demodulation. Blue and red lines show their modulators obtained by using the AP-A algorithm with the cutoff frequency $\omega$ set to, respectively, 55 Hz and 222 Hz. Insets display time-expanded segments of the original window. Red arrowheads indicate the ringing artifacts of the modulator at some prolonged intervals of low signal levels. Source of the original signal: audio edition of *The Economist* magazine, issue March 19th 2016, article "Restoring lost memories."

of $\hat{\mathbf{m}}$ and $|\mathbf{s}|$ at these sample points. We illustrate our claims in Fig. 6 A, where a band-pass component of a female utterance *"...at one hundred hertz ..."* with $f_c = 800$ Hz, $\Delta = 111$ Hz, and $\varpi = 222$ Hz is considered. Taking into account that the local maximum points $|s_i|$ are locally regular and that they correspond to $m_i$, we could exploit *Proposition A.2* in Suppl. Mat. A to find that $E_m \leq 8 \cdot 10^{-3}$.

Wideband speech signals are more challenging than their narrow subbands. They are built of temporarily structured segments of quasi-random and quasi-harmonic carriers, possibly featuring frequency glides [65]. These carriers are amplitude-modulated at different timescales, ranging between a hundred milliseconds and several seconds [23], [66]. The power spectral density of the corresponding modulators is vanishingly small above 20 Hz (see Fig. 1 in [67]). Moreover, as we demonstrate in Suppl. Mat. O, the carrier components of natural speech signals align to the recoverability conditions of the AP approach for $\mathbf{m} \in \mathcal{M}_\omega$ with $\omega$ up to at least $\sim 50$ Hz. Therefore, we expect appropriate performance from the AP algorithms in the setting of wideband speech.

Fig. 6 B displays demodulation results of the full-band version of the speech segment considered in Fig. 6 A by the AP-A algorithm with $\varpi = 55$ Hz. The obtained $\hat{\mathbf{m}}$ (red) envelops separate phonemes of the sound waveform tightly, indicating appropriate recovery of the true $\mathbf{m}$ (see Section VIII-B next). However, intervals corresponding to prolonged transitions between phonemes or words are corrupted by ringing artifacts (marked by red arrowheads in Fig. 6 B), implying the necessity of higher frequency components to represent these transitions. Hence, although the power spectral density of the true $\mathbf{m}$ is very low above 20 Hz, it sums to a noticeable contribution. Unfortunately, any attempt to cancel the artifacts by just increasing $\varpi$ fails by breaking the recovery conditions, as

illustrated by the blue line in Fig. 6 B ($\varpi = 222$ Hz there). No improvement is achieved by utilizing the AP-B, AP-P, or LDC algorithms either (data not shown).

### B. Demodulation using dynamic range compression

The aforementioned problem with modulator estimates of signals with sharp transitions to/from prolonged intervals of low-signal amplitude can be resolved by using a dynamic range compression. In particular, instead of demodulating the original signal $\mathbf{s}$ directly, we first apply a chosen AP algorithm to its compressed version:

$$\underline{\mathbf{s}} = \text{sgn}(\mathbf{s}) \circ |\mathbf{s}|^{1/p}. \tag{12}$$

Here, $p \in (1, +\infty)$ controls the level of compression. The modulator estimate $\hat{\mathbf{m}}^*$ of $\mathbf{s}$ is then evaluated by inverse-transforming the modulator $\underline{\hat{\mathbf{m}}}$ of $\underline{\mathbf{s}}$:

$$\hat{\mathbf{m}}^* = \underline{\hat{\mathbf{m}}}^p. \tag{13}$$

The idea behind (12) is that the compression makes signals more uniform and, effectively, smooths their sharp changes responsible for ringing artifacts in the modulator estimates. These sharp changes are restored in the modulators without artifacts by the inverse transform (13).

The expected effect of the compression procedure is illustrated in Fig. 7, where signal demodulation of an utterance *"...protein which forms p..."* is considered. Differently from the direct demodulation result $\hat{\mathbf{m}}$ (red line), the estimate $\hat{\mathbf{m}}^*$ obtained by using the compression with $p = 3$ (black line) shows good alignment with $|\mathbf{s}|$ in the segments of both low and high intensity. To justify that this alignment really reflects the recovery of the true modulator, we performed additional tests where chimeric signals built of $\hat{\mathbf{m}}^*$ from Fig. 7 and natural speech carriers were demodulated (see Suppl. Mat. O).
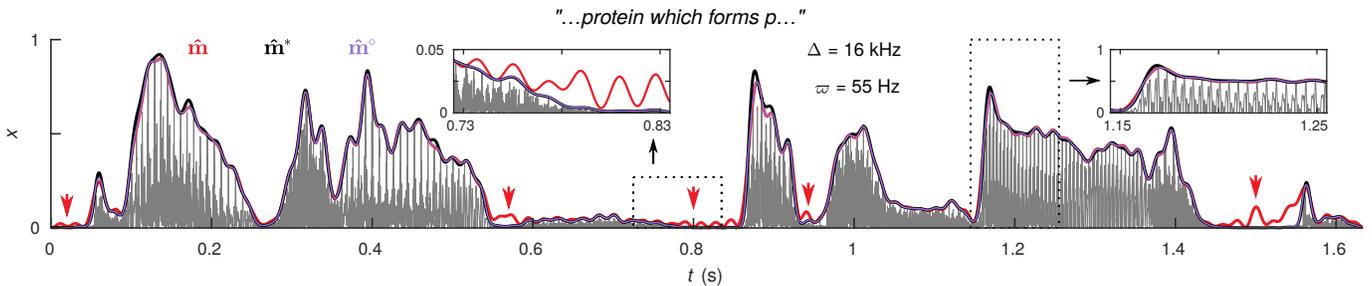
Fig. 7. Demodulation of speech signals using dynamic range compression. An audio signal of *"...protein which forms p ..."* uttered by a female speaker (full bandwidth of 16 kHz). Gray – the absolute-value version of the signal considered for demodulation, red – its modulator obtained by using the AP-A algorithm with $\omega = 55$ Hz and no compression (as in Fig. 6), black – modulator of the signal obtained when employing the dynamic range compression [see (13)], violet – interpolation of the latter two [see (14)]. Red arrowheads indicate ringing artifacts of the modulator estimate. Source of the original signal: the same as Fig. 6.

We found low demodulation errors, with $E_m$ ranging between $9 \cdot 10^{-3}$ and $5 \cdot 10^{-2}$ for different carrier components of speech signals (see Fig. 18).

The compression level $p = 3$ used above was adjusted by a trial and error for speech signals. In general, the gains in accuracy at low levels with increased $p$ comes at the expense of reduced precision of modulator estimated at high signal levels. Thus, a compromise between those two effects must be reached to find an optimal $p$. Moreover, the precision of the modulator estimates can be further increased by interpolating between $\hat{\mathbf{m}}^*$ (more accurate for low signal levels) and $\hat{\mathbf{m}}$ (more accurate for high signal levels). For example, the violet line in Fig. 7 shows a weighted average of the form

$$\hat{\mathbf{m}}^{\diamond} = \hat{\mathbf{m}} \circ \mathbf{w} + \hat{\mathbf{m}}^* \circ (1 - \mathbf{w}), \qquad (14)$$

where

$$w_i = \left( \frac{1 - e^{a \cdot (\hat{m}_i^* / \max[\hat{\mathbf{m}}^*])}}{1 + e^{a \cdot (\hat{m}_i^* / \max[\hat{\mathbf{m}}^*]) - b}} \right) \cdot \left( \frac{1 - e^a}{1 + e^{a-b}} \right)^{-1} \qquad (15)$$

for $i \in \mathcal{I}_n$, with $b = 3$ and $a = 10$. In general, an optimal interpolation between $\hat{\mathbf{m}}$ and $\hat{\mathbf{m}}^*$ can be learned by minimizing $\|\hat{\mathbf{m}}^{\diamond}\|_2^2$ over a chosen class of functions. Other compression models than (12), e.g., $\underline{\mathbf{s}} = \mathrm{sgn}(\mathbf{s}) \circ \log(1 + p \cdot |\mathbf{s}|)$, can be used to evaluate $\hat{\mathbf{m}}^*$ as well.

### C. Demodulation in real-time

A number of amplitude demodulation applications, e.g., speech recognition [2], ultrasound imaging [68], and cochlear prosthesis [55], necessitate real-time processing. As we demonstrate below, the exceptional computational efficiency of the AP approach allows it to fulfill that requirement.

The nature of the task implies that online modulator estimates have to be generated by sequentially demodulating windowed segments $\mathbf{s}^{(j)}$ of a signal $\mathbf{s}$ at each updated sample point $j$ across time:

$$\mathbf{s}^{(j)} : s_i^{(j)} = w_i \cdot s_{j-k_l-1+i}, \qquad i \in \{1, 2, \dots, k\}. \quad (16)$$

Here, $k$ is the number of sample points corresponding to the segment, and $k_l$ denotes the number of sample points of it that are to the left of the current point $j$. $w_i, w_2, \dots, w_k$ are vector

elements of the window function. The real-time modulator estimate $\hat{m}_j^{\star}$ at sample point $j$ is calculated as

$$\hat{m}_j^{\star} = \hat{m}_{k_l+1}^{(j)}, \qquad (17)$$

where $\hat{\mathbf{m}}^{(j)}$ is a modulator estimate of $\mathbf{s}^{(j)}$.

It follows from the time-frequency uncertainty principle [8] that accurate evaluation of $\hat{m}_j^{\star}$ requires $\mathbf{s}^{(j)}$ with a duration of the order of the inverse of the effective bandwidth of the modulator, or longer. This condition sets the lower bounds on the segment length $k$ and sampling delay $k_\tau = k - k_l - 1$ of $\hat{\mathbf{m}}^{\star}$. We found empirically that $k_\tau \approx 2 \cdot (f_s / \varpi)$ and $k \approx 4 \cdot (f_s / \varpi)$ are typically sufficient for accurate demodulation of wideband speech. These numbers are around two times smaller for narrow frequency band components of these signals. We know that $\varpi \geq 40$ Hz for the wideband speech and its subbands. Thus, delays $k_\tau \leq 50$ ms for estimating $\hat{\mathbf{m}}^{\star}$ are sufficient without a sacrifice in precision then. The main requirement for the window function in (16) is that it smoothly scales the signal to 0 at the boundaries, with no effect at the midst. We used a modified version of the Hann window for this purpose:

$$w_i = \begin{cases} \sin^2 \left( \frac{\pi \cdot (i-1)}{2 \cdot k_l} \right), & 1 \leq i \leq k_l \\ 1, & i = k_l + 1 \\ \cos^2 \left( \frac{\pi \cdot (i-k+k_\tau)}{2 \cdot k_\tau} \right), & k - k_\tau + 1 \leq i \leq k \end{cases} . \quad (18)$$

Fig. 8 shows simulation results of real-time demodulation of a male utterance *"...with little human hand-holding ..."* (sampling rate $f_s = 16$ kHz) based on the AP-A algorithm. There, demodulation was performed with $k = 1536$ and $k_\tau = 768$ ($\tau = 48$ ms) for the original signal (Fig. 8 B). Its subband component centered at 800 Hz (Fig. 8 A) was processed with $k = 128$ and $k_\tau = 65$ ($\tau = 4$ ms). In each case, $\hat{m}_j^{\star}$ was updated with the frequency of $10 \cdot \varpi$. The obtained estimates $\hat{\mathbf{m}}^{\star}$ are in very good agreement with $\hat{\mathbf{m}}^*$ derived by using offline demodulation of the whole signal, with $\|\hat{\mathbf{m}}^{\star} - \hat{\mathbf{m}}^{\diamond}\|_2 / \|\hat{\mathbf{m}}^{\diamond}\|_2 < 0.02$. Importantly, they were achieved with modest CPU usage: $T_{\mathrm{cpu}}$ amounted to only 1.6 % (subband signal) and 3.2 % (wideband signal) of the time length of the demodulated signal on an Intel Core i7-7700 CPU run in single-thread mode. For comparison, these
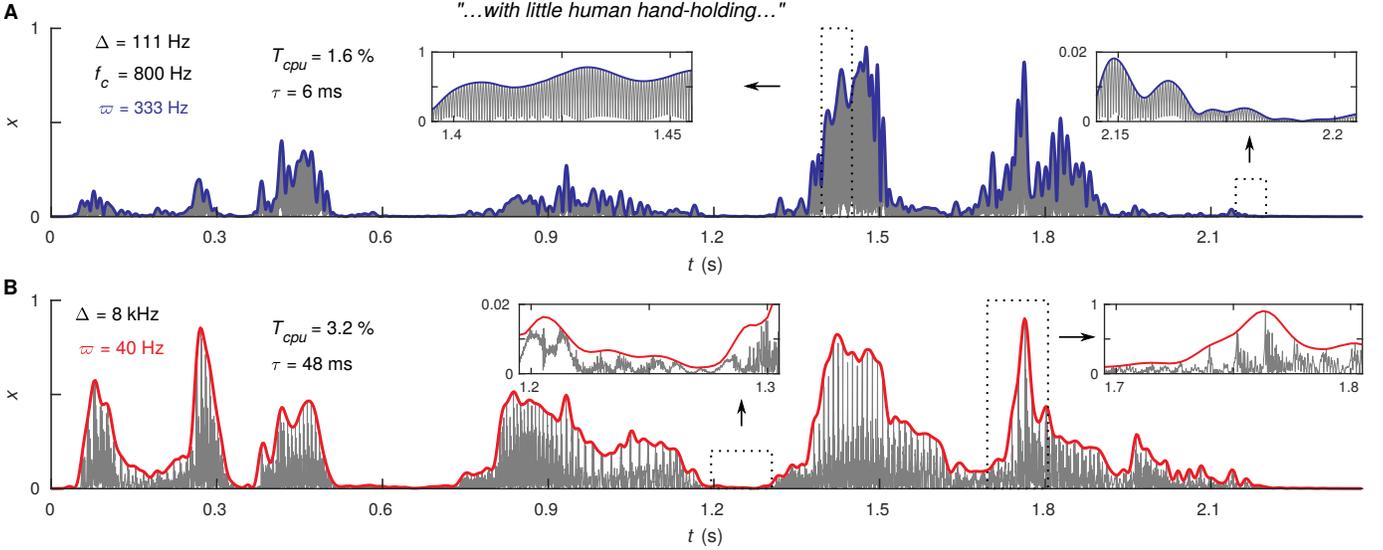
Fig. 8. Demodulation in real-time. **A**: A band-pass-filtered signal of an utterance *"…with little human hand-holding …"* by a male speaker; the signal was obtained with an equivalent rectangular bandwidth filter of the cochlea centered at 800 Hz (bandwidth of 111 Hz). **B**: The original signal of the utterance used in panel A (full bandwidth of 8 kHz). In both panels, gray color marks the absolute-value version of the signals considered for demodulation. Blue and red lines show their modulators obtained by using the real-time version of the AP-A algorithm with the cutoff frequency $\omega$ set to, respectively, 333 Hz and 40 Hz. Source of the original signal: audio edition of *The Economist* magazine, issue March 19th 2016, article "Artificial intelligence and Go."

numbers were, respectively, $\sim 5 \cdot 10^3$ and $\sim 6 \cdot 10^4$ times higher for the LDC method.

An advantageous side effect of splitting the signal into small windows for demodulation is that it prevents the ringing artifacts (compare Fig. 8 B and Fig. 6 B). This is so because signal levels do not typically spread over different scales in a short time window. The window splitting also allows generalizing demodulation to situations when the cutoff frequency $\omega$ of the modulator varies strongly in time.

## IX. EXTENSIONS AND GENERALIZATIONS

### A. Demodulation in higher dimensions

Amplitude demodulation has found successful applications beyond the setting of 1D signals. Several 2D extensions of the classical AS approach have been introduced and used for solving tasks in computer vision [7], [69], analysis of speech spectrograms [70], and biomedical imaging [4], [71], [72]. The AS framework has also been extended to calculate modulators and carriers for signals over graphs [73]. These methods are limited to locally narrowband signals, which manifest visually as fringe patterns (see Fig. 9 A). This bandwidth restriction is evaded by a generalization of the AP approach to higher dimensions that we present next. The extension is immediate and follows from intuitive abstractions of the constraint sets introduced in Section II.

Consider a $D$-dimensional signal $s(t_1, t_2, \ldots, t_D)$. Its uniformly sampled version $\bar{\mathbf{s}}$ is an element of an $n$-dimensional Euclidean space $\mathbb{T}_D^n$ of real-valued order $D$ tensors with $n = \prod_{i=1}^{D} n_i$ and the inner product $\langle \bar{\mathbf{s}}^{(1)}, \bar{\mathbf{s}}^{(2)} \rangle = \sum_{i_1=1}^{n_1} \cdots \sum_{i_D=1}^{n_D} (\bar{s}_{i_1 \cdots i_D}^{(1)} \cdot \bar{s}_{i_1 \cdots i_D}^{(2)})$. The respective $D$-dimensional DFT is given by

$$\overline{\mathbf{F}} = \mathbf{F}^{(1)} \otimes \mathbf{F}^{(2)} \otimes \cdots \otimes \mathbf{F}^{(D)}, \qquad (19)$$

where $\mathbf{F}^{(i)}$ is a unitary DFT defined over $\mathbb{R}^{n_i}$. Then, the analogs of the constraint sets $\mathcal{S}_{\geq \mathbf{0}}$, $\mathcal{S}_\omega$, $\mathcal{S}_{\geq |\mathbf{s}|}$, $\mathcal{S}_{|\cdot| \leq \mathbf{1}}$, and $\mathcal{S}_{\{1\}, d}$ from Section II read as

$$\mathcal{S}_{\geq \bar{\mathbf{0}}} = \{\bar{\mathbf{x}} \in \mathbb{T}_D^n : \bar{x}_{i_1 \cdots i_D} \geq 0, \ i_j \in \mathcal{I}_{n_j}\},$$
$$\mathcal{S}_\omega = \{\bar{\mathbf{x}} \in \mathbb{T}_D^n : (\overline{\mathbf{F}}\,\bar{\mathbf{x}})_{i_1 \cdots i_D} = 0, \ i_j \in (\mathcal{I}_{n_j} \setminus \mathcal{I}_{n_j}^{\omega_j})\}, \quad (20)$$
$$\mathcal{S}_{\geq |\bar{\mathbf{s}}|} = \{\bar{\mathbf{x}} \in \mathbb{T}_D^n : \bar{x}_{i_1 \cdots i_D} \geq |\bar{s}_{i_1 \cdots i_D}|, \ i_j \in \mathcal{I}_{n_j}\},$$

and

$$\mathcal{S}_{|\cdot| \leq \bar{\mathbf{1}}} = \{\bar{\mathbf{x}} \in \mathbb{T}_D^n : |x_{i_1 \cdots i_D}| \leq 1, \ i_j \in \mathcal{I}_{n_j}\},$$
$$\mathcal{S}_{\{1\}, \mathbf{d}} = \{\bar{\mathbf{x}} \in \mathbb{T}_D^n : (\forall i_1 \cdots i_D) R(i_1 \cdots i_D, \bar{\mathbf{x}}, \mathbf{d}) \geq 1, \quad (21)$$
$$(\exists i_1 \cdots i_D) R(i_1 \cdots i_D, \bar{\mathbf{x}}, \mathbf{d}) = 1\},$$

where

$$R(i_1 \cdots i_D, \bar{\mathbf{x}}, \mathbf{d}) = \sum_{j_1 \geq i_1} \cdots \sum_{j_D \geq i_D} \left[ I_{\{1\}}(|\bar{x}_{j_1 \cdots j_D}|) \right.$$
$$\left. \cdot \theta\big(1 - \sum_{k=1}^{D}(i_k - j_k)^2/d_k^2\big) \right] - I_{\{1\}}(|\bar{x}_{i_1 \cdots i_D}|). \quad (22)$$

Simply substituting (20)–(21) for their $D = 1$ versions in (2), (4), (6), (8), (9), and (10) generalizes the modulator $\mathcal{M}_\omega$ and carrier $\mathcal{C}_d$ sets as well as the modulator estimator $\hat{\mathbf{m}}$ and the respective AP algorithms. In particular, an $\overline{\mathbf{m}} \in \mathcal{M}_\omega$ is a nonnegative signal with a low-pass rectangular spectrum set by $\boldsymbol{\omega} = (\omega_1, \ldots, \omega_D)$ along each of the $D$ dimensions in the DFT domain. A $\bar{\mathbf{c}} \in \mathcal{C}_\mathbf{d}$ is a signal bounded between $-1$ and 1 with the $|\bar{c}_{i_1 \cdots i_D}| = 1$ sample points packed sufficiently densely, as implied by $\mathbf{d} = (d_1, \ldots, d_D)$.

Without providing formal proofs, we state that all propositions and assertions of Sections II and III about the modulator recoverability and convergence of the AP algorithms generalize to $D$-dimensional signals defined above. All quantitative conditions involving the parameters $\varpi$, $\omega$, $d$, $n$, and $n_s$ in the $D = 1$ case are then replaced by elementwise conditions for $\varpi_i$, $\omega_i$, $d_i$, $n_i$, and $n_{s,i}$ at $i \in \mathcal{I}_D$.
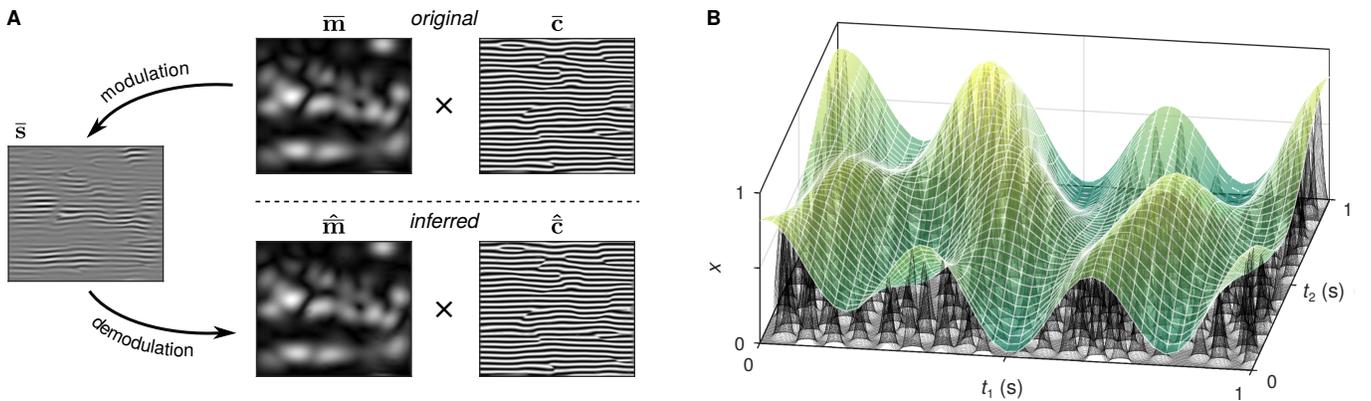
Fig. 9. Demodulation in 2D by using the AP-A algorithm. **A**: Synthetic fringe pattern of moderate bandwidth. Left – the pattern to demodulate, upper right – the original modulator and carrier, lower right – the inferred modulator and carrier. **B**: Wideband signal consisting of randomly-placed spikes of finite width. Black grid – the signal $\bar{\mathbf{s}}$ to demodulate, white grid – the original modulator $\overline{\mathbf{m}}$ of the signal, color surface – the estimated modulator.

Fig. 9 illustrates the potential of the AP-A algorithm with the help of two $D = 2$ cases. Fig. 9 A shows successful demodulation results for a synthetic narrowband fringe pattern ($E_m = 1 \cdot 10^{-2}$, $E_c = 3 \cdot 10^{-2}$). Fig. 9 B displays high-accuracy demodulation of a wideband signal built of randomly-placed spikes of finite width as $\overline{\mathbf{c}}$ and a Gaussian random field with a rectangular amplitude spectrum as $\overline{\mathbf{m}}$. There, the white grid corresponds to the original $\overline{\mathbf{m}}$, while the color surface represents its estimate $\hat{\overline{\mathbf{m}}}$ ($E_m = 4 \cdot 10^{-3}$, $E_c = 1 \cdot 10^{-2}$).

The ability of the AP approach to deal with wideband signals allows it to cover a wider range of practically relevant situations. Among examples are nonlinear ultrasound imaging [4], [14], speech processing [20], [70], and complicated cases of optical interference/diffraction setups [74]. Moreover, it can also be of great use in time-critical imaging settings by providing high modulator estimation accuracy at low sampling rates of the signal (see, e.g., [72], [75]).

The minimum number of sample points necessary to cover simultaneously for appropriate demodulation increases exponentially with $D$. Therefore, the computational advantage of the AP over the PAD and LDC demodulation approaches is even more pronounced in higher dimensions. In fact, if evaluated by using the FFT method, $\overline{\mathbf{F}}$ features an $\mathcal{O}(n \log n)$ computational time complexity. Hence, the time complexity of the AP algorithms is defined by the total number of sample points of the signal irrespective of its dimensionality.

### B. Generalized modulators and nonuniform sampling

The demodulation approach formulated in the present work builds on the assumption that modulators are nonnegative elements of a low-pass DFT subspace of $\mathbb{T}_D^n$. However, as follows from the convergence proofs in Suppl. Mat. F, all of the introduced AP algorithms are bound to converge to an $\hat{\mathbf{m}} \in \mathcal{M}_\omega$ and a $\hat{\mathbf{c}} \in \mathcal{C}_{\mathbf{d}}$ independent of the origin of the linear subspace behind $\mathcal{M}_\omega$. This naturally raises the question of whether the AP algorithms could recover true $\mathbf{m}$ and $\mathbf{c}$ under the generalized subspace assumption. Our preliminary experiments suggest a positive answer but subject to extra recovery conditions specific to a subspace of choice.

For example, consider a subset of $2\omega - 1$ randomly chosen basis vectors of the DFT over $\mathbb{R}^n$. Denote the corresponding space as $\mathcal{F}_\omega$. It can be shown by example that a system resulting from random subsampling of the aforementioned vectors at $2\omega - 1$ time points may be linearly dependent. If so, it then follows from the proof of *Proposition II.1* that, in contrast to an $\mathbf{m} \in \mathcal{S}_\omega$, full recovery of an $\mathbf{m} \in \mathcal{F}_\omega$ necessitates more than $2\omega - 1$ true sample points.

The problem of formulating modulator recovery conditions for different linear subspaces sets directions for future studies. If successful, these extensions would allow to:

1) broaden the concept of the amplitude modulator beyond the low-pass DFT signals,
2) loosen the constraints on the positioning of the $|c_i| = 1$ sample points for recoverable carriers whenever a more compact representation of modulators is available,
3) encompass nonuniform sampling.

While the above points are yet to be developed, the results of the present work already provide a strategy for an arbitrarily-accurate nonuniform sampling. Indeed, for any time grid $\tilde{\mathbf{t}} \in \mathbb{R}^{\tilde{n}}$, we can find a uniform grid $\mathbf{t} \in \mathbb{R}^n$ such that, for every $j \in \mathcal{I}_{\tilde{n}}$, there exists an $i \in \mathcal{I}_n$ with $|\tilde{t}_j - t_i|$ being arbitrarily small. We can then interpolate the original data $\tilde{\mathbf{s}} \in \mathbb{R}^{\tilde{n}}$ on the uniform grid $\mathbf{t}$ by

$$s_i = \begin{cases} \tilde{s}_j, & \text{if } |\tilde{t}_j - t_i| = \min[|\tilde{t}_j \cdot \mathbf{1} - \mathbf{t}|] \\ 0, & \text{otherwise} \end{cases}, \quad i \in \mathcal{I}_n, \quad (23)$$

to obtain an $\mathbf{s} \in \mathbb{R}^n$. The bandwidth constraint on $\mathbf{m}$ implies that all components of $\mathbf{s}$ corresponding to the true sample points of $\tilde{\mathbf{s}}$ are desirably close to the true sample points of the original signal if $n$ is large enough. Then, *Proposition II.2* assures that modulator-carrier recovery is possible via (6) under the conditions discussed in Section II-C for uniformly sampled signals. The described strategy requires increasing the effective dimensionality of the signal. However, this may still be more efficient than evaluating metric projections onto subspaces spanned by arbitrary nonuniform sampling basis vectors, which are not orthogonal in general.

## X. Conclusion

In this paper, we have introduced a new approach to amplitude demodulation of arbitrary-bandwidth signals. We framed demodulation as a problem of modulator recovery from an unlabeled mix of its true and corrupted sample points. Taking this view, we showed that high-accuracy demodulation can be achieved via exact or approximate norm minimization of the modulator for a wide range of relevant signals. We formulated tailor-made alternating projection algorithms to achieve that in practice and tested them in a series of numerical experiments.

The generality and numerical efficiency of the new approach make it a preferred choice in many situations. In the context of narrowband signals, the new method outperforms the classical algorithms in terms of robustness to data distortions and compatibility with nonuniform sampling. When considering the demodulation of wideband signals, it surpasses the current state-of-the-art techniques in terms of computational efficiency by up to many orders of magnitude. Such performance enables practical applications of amplitude demodulation in previously inaccessible settings. Specifically, online and large-scale offline demodulation of wideband signals, signals in higher dimensions, and poorly-sampled signals become practically feasible. The algorithms underlying the new approach are simple and easy to implement on a computer.[5]

## Acknowledgment

## References

[1] D. Vakman, *Signals, oscillations, and waves: A modern approach.* Artech House, 1998.

[2] B. E. D. Kingsbury, N. Morgan, and S. Greenberg, "Robust speech recognition using the modulation spectrogram," *Speech Commun.*, vol. 25, pp. 117–132, 1998.

[3] M. G. Ruppert, D. M. Harcombe, M. R. P. Ragazzon, S. O. R. Moheimani, and A. J. Fleming, "A review of demodulation techniques for amplitude-modulation atomic force microscopy," *Beilstein J. Nanotechnol.*, vol. 8, pp. 1407–1426, 2017.

[4] C. Wachinger, T. Klein, and N. Navab, "The 2D analytic signal for envelope detection and feature extraction on ultrasound images," *Medical Image Analysis*, vol. 16, pp. 1073–1084, 2012.

[5] P. Y. Ktonas and N. Papp, "Instantaneous envelope and phase extraction from real signals: Theory, implementation, and an application to EEG analysis," *Elsevier Signal Process.*, vol. 2, pp. 373–385, 1980.

[6] M. T. Taner, F. Koehler, and R. E. Sheriff, "Complex seismic trace analysis," *Geophysics*, vol. 44, pp. 1041–1063, 1979.

[7] K. G. Larkin, D. J. Bone, and M. A. Oldfield, "Natural demodulation of two-dimensional fringe patterns. I." *J. Opt. Soc. Am. A*, vol. 18, pp. 1862–1870, 2001.

[8] D. Gabor, "Theory of communication. Part 1: The analysis of information," *J. Inst. Elec. Eng. Part III*, vol. 93, pp. 429–441, 1946.

[9] D. Vakman, "On the analytic signal, the Teager-Kaiser energy algorithm, and other methods for defining amplitude and frequency," *IEEE Trans. Signal Process.*, vol. 44, pp. 791–797, 1996.

[10] B. S. Wilson, C. C. Finley, D. T. Lawson, R. D. Wolford, D. K. Eddington, and W. M. Rabinowitz, "Better speech recognition with cochlear implants," *Nature*, vol. 352, pp. 236–238, 1991.

[11] Z. M. Smith, B. Delgutte, and A. J. Oxenham, "Chimaeric sounds reveal dichotomies in auditory perception," *Nature*, vol. 416, pp. 87–90, 2002.

[12] U. Goswami, "Speech rhythm and language acquisition: An amplitude modulation phase hierarchy perspective," *Ann. N. Y. Acad. Sci.*, vol. 1453, pp. 67–78, 2019.

[13] S. Lin, "Demodulating wide-band ultrasound signals," U.S. Patent US6 248 071B1, 2001.

[14] F. A. Duck, "Nonlinear acoustics in diagnostic ultrasound," *Ultrasound Med. Biol.*, vol. 28, pp. 1–18, 2002.

[15] G. L. Gottlieb and G. C. Agarwal, "Filtering of electromyographic signals," *Am. J. Phys. Med. Rehab.*, vol. 49, p. 142, 1970.

[16] J. Felblinger and C. Boesch, "Amplitude demodulation of the electrocardiogram signal (ECG) for respiration monitoring and compensation during MR examinations," *Magn. Reson. Med.*, vol. 38, pp. 129–136, 1997.

[17] D. Gill, N. Gavrieli, and N. Intrator, "Detection and identification of heart sounds using homomorphic envelogram and self-organizing probabilistic model," in *Computers in Cardiology, 2005*, 2005, pp. 957–960.

[18] W. Liu and B. Santhanam, "Wideband image demodulation via bidimensional multirate frequency transformations," *J. Opt. Soc. Am. A*, vol. 33, pp. 1668–1678, 2016.

[19] G. Sell and M. Slaney, "Solving demodulation as an optimization problem," *IEEE Audio, Speech, Language Process.*, vol. 18, pp. 2051–2066, 2010.

[20] ——, "The information content of demodulated speech," in *IEEE Proc. ICASSP'10*, 2010, pp. 5470–5473.

[21] R. Libbey, *Signal & image processing sourcebook.* Springer, 1994.

[22] R. S. Platt, E. A. Hajduk, M. Hulliger, and P. A. Easton, "A modified bessel filter for amplitude demodulation of respiratory electromyograms," *J. Appl. Physiol.*, vol. 84, pp. 378–388, 1998.

[23] R. E. Turner, "Statistical models for natural sounds," Ph.D. dissertation, University College London, 2010. [Online]. Available: http://discovery.ucl.ac.uk/19231/

[24] R. E. Turner and M. Sahani, "Demodulation as probabilistic inference," *IEEE Audio, Speech, Language Process.*, vol. 19, pp. 2398–2411, 2011.

[25] P. J. Loughlin and B. Tacer, "On the amplitude- and frequency-modulation decomposition of signals," *J. Acoust. Soc. Am.*, vol. 100, pp. 1594–1601, 1996.

[26] L. Cohen, P. Loughlin, and D. Vakman, "On an ambiguity in the definition of the amplitude and phase of a signal," *Elsevier Signal Process.*, vol. 79, pp. 301–307, 1999.

[27] J. von Neumann, *Functional operators. Vol. II: The geometry of orthogonal spaces*, ser. Annals of Mathematics Studies 22. Princeton University Press, 1951.

[28] R. Escalante and M. Raydan, *Alternating projection methods.* SIAM, 2011.

[29] H. H. Bauschke and J. M. Borwein, "On projection algorithms for solving convex feasibility problems," *SIAM Rev.*, vol. 38, pp. 367–426, 1996.

[30] A. Ahmed, B. Recht, and J. Romberg, "Blind deconvolution using convex programming," *IEEE Trans. Inf. Theory*, vol. 60, pp. 1711–1732, 2014.

[31] S. Ling and T. Strohmer, "Self-calibration and biconvex compressive sensing," *Inverse Probl.*, vol. 31, p. 115002, 2015.

[32] Y. Chi, "Guaranteed blind sparse spikes deconvolution via lifting and convex optimization," *IEEE J. Sel. Topics Signal Process.*, vol. 10, pp. 782–794, 2016.

[33] Y. Xie, M. B. Wakin, and G. Tang, "Simultaneous sparse recovery and blind demodulation," *IEEE Trans. Signal Process.*, vol. 67, pp. 5184–5199, 2019.

[34] A. Oppenheim and J. Lim, "The importance of phase in signals," *Proc. IEEE*, vol. 69, pp. 529–541, 1981.

[35] H. Trussell and M. Civanlar, "The feasible solution in signal restoration," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 32, pp. 201–212, 1984.

[36] D. Kundur and D. Hatzinakos, "A novel blind deconvolution scheme for image restoration using recursive filtering," *IEEE Trans. Signal Process.*, vol. 46, pp. 375–390, 1998.

[37] Y. Yang, N. P. Galatsanos, and H. Stark, "Projection-based blind deconvolution," *J. Opt. Soc. Am. A*, vol. 11, pp. 2401–2409, 1994.

[38] P. J. S. G. Ferreira, "Iterative and noniterative recovery of missing samples for 1-D band-limited signals," in *Nonuniform sampling: Theory and practice*, F. Marvasti, Ed. Springer, 2001, pp. 235–281.

[39] L. G. Gubin, B. T. Polyak, and E. V. Raik, "The method of projections for finding the common point of convex sets," *USSR Comput. Math. & Math. Phys.*, vol. 7, pp. 1–24, 1967.

[40] D. C. Youla and H. Webb, "Image restoration by the method of convex projections: Part 1– theory," *IEEE Trans. Med. Imag.*, vol. 1, pp. 81–94, 1982.

[41] C. Franchetti and W. Light, "On the von Neumann alternating algorithm in Hilbert space," *J. Math. Anal. Appl.*, vol. 114, pp. 305–314, 1986.

---

[5]The computer code for AP demodulation will be available at *https://github.com/mgabriel-lt/ap-demodulation*.

[42] W. B. Gearhart and M. Koshy, "Acceleration schemes for the method of alternating projections," *J. Comput. Appl. Math.*, vol. 26, pp. 235–249, 1989.

[43] H. H. Bauschke, F. Deutsch, H. Hundal, and S.-H. Park, "Accelerating the convergence of the method of alternating projections," *Trans. Amer. Math. Soc.*, vol. 355, pp. 3433–3461, 2003.

[44] R. L. Dykstra, "An algorithm for restricted least squares regression," *J. Amer. Statist. Assoc.*, vol. 78, pp. 837–842, 1983.

[45] J. P. Boyle and R. L. Dykstra, "A method for finding projections onto the intersection of convex sets in Hilbert spaces," in *Advances in Order Restricted Statistical Inference*, ser. Lecture Notes in Statistics. Springer, 1986, pp. 28–47.

[46] E. Birgin and M. Raydan, "Robust stopping criteria for Dykstra's algorithm," *SIAM J. Sci. Comput.*, vol. 26, pp. 1405–1414, 2005.

[47] F. Deutsch and H. Hundal, "The rate of convergence of Dykstra's cyclic projections algorithm: The polyhedral case," *Numer. Funct. Anal. Optim.*, vol. 15, pp. 537–565, 1994.

[48] F. Deutsch, "Dykstra's cyclic projections algorithm: The rate of convergence," in *Approximation Theory, Wavelets and Applications*, ser. NATO Science Series. Springer, 1995, pp. 87–94.

[49] H. H. Bauschke and J. M. Borwein, "Dykstra's alternating projection algorithm for two sets," *J. Approx. Theory*, vol. 79, pp. 418–443, 1994.

[50] P. Duhamel and M. Vetterli, "Fast Fourier transforms: A tutorial review and a state of the art," *Elsevier Signal Process.*, vol. 19, pp. 259–299, 1990.

[51] Gurobi Optimization, LLC, *Gurobi optimizer reference manual*, 2019. [Online]. Available: http://www.gurobi.com

[52] B. Stellato, G. Banjac, P. Goulart, A. Bemporad, and S. Boyd, "OSQP: an operator splitting solver for quadratic programs," *Math. Prog. Comp.*, vol. 12, pp. 637–672, 2020.

[53] L. Marple, "Computing the discrete-time "analytic" signal via FFT," *IEEE Trans. Signal Process.*, vol. 47, pp. 2600–2603, 1999.

[54] R. Wiley, H. Schwarzlander, and D. Weiner, "Demodulation Procedure for Very Wide-Band FM," *IEEE Trans. Commun.*, vol. 25, pp. 318–327, 1977.

[55] B. S. Wilson and M. F. Dorman, "Cochlear implants: A remarkable past and a brilliant future," *Hear. Res.*, vol. 242, pp. 3–21, 2008.

[56] S. Wu, T. H. Falk, and W.-Y. Chan, "Automatic speech emotion recognition using modulation spectral features," *Speech Commun.*, vol. 53, pp. 768–785, 2011.

[57] B. Lee and K.-H. Cho, "Brain-inspired speech segmentation for automatic speech recognition using the speech envelope as a temporal reference," *Sci. Rep.*, vol. 6, p. 37647, 2016.

[58] G. Hu and D. Wang, "Monaural speech segregation based on pitch tracking and amplitude modulation," *IEEE Trans. Neural Netw.*, vol. 15, pp. 1135–1150, 2004.

[59] L. Atlas and C. Janssen, "Coherent modulation spectral filtering for single-channel music source separation," in *IEEE Proc. ICASSP'05*, vol. 4, 2005, pp. 461–464.

[60] P. X. Joris, C. E. Schreiner, and A. Rees, "Neural processing of amplitude-modulated sounds," *Physiol. Rev.*, vol. 84, pp. 541–577, 2004.

[61] F.-G. Zeng, K. Nie, G. S. Stickney, Y.-Y. Kong, M. Vongphoe, A. Bhargave, C. Wei, and K. Cao, "Speech recognition with amplitude and frequency modulations," *PNAS*, vol. 102, pp. 2293–2298, 2005.

[62] R. V. Shannon, F.-G. Zeng, V. Kamath, J. Wygonski, and M. Ekelid, "Speech recognition with primarily temporal cues," *Science*, vol. 270, pp. 303–304, 1995.

[63] B. R. Glasberg and B. C. J. Moore, "Derivation of auditory filter shapes from notched-noise data," *Hear. Res.*, vol. 47, pp. 103–138, 1990.

[64] J. L. Flanagan, "Parametric coding of speech spectra," *J. Acoust. Soc. Am.*, vol. 68, pp. 412–419, 1980.

[65] J. E. Shoup and L. L. Pfeifer, "Acoustic characteristics of speech sounds," in *Contemporary Issues in Experimental Phonetics*. Academic Press, 1976, pp. 171–224.

[66] A. Keitel, J. Gross, and C. Kayser, "Perceptually relevant speech tracking in auditory and motor cortex reflects distinct linguistic features," *PLOS Biol.*, vol. 16, p. e2004473, 2018.

[67] H. R. Bosker and M. Cooke, "Talkers produce more pronounced amplitude modulations when speaking in noise," *J. Acoust. Soc. Am.*, vol. 143, pp. EL121–EL126, 2018.

[68] P. R. Hoskins, K. Martin, and A. Thrush, *Diagnostic ultrasound: Physics and equipment*, 3rd ed. CRC Press, 2019.

[69] M. Felsberg and G. Sommer, "The monogenic signal," *IEEE Trans. Signal Process.*, vol. 49, pp. 3136–3144, 2001.

[70] H. Aragonda and C. S. Seelamantula, "Demodulation of narrowband speech spectrograms using the Riesz transform," *IEEE Audio, Speech, Language Process.*, vol. 23, pp. 1824–1834, 2015.

[71] C. S. Seelamantula, N. Pavillon, C. Depeursinge, and M. Unser, "Local demodulation of holograms using the Riesz transform with application to microscopy," *J. Opt. Soc. Am. A*, vol. 29, pp. 2118–2129, 2012.

[72] K. Nadeau, A. J. Durkin, and B. J. Tromberg, "Advanced demodulation technique for the extraction of tissue optical properties and structural orientation contrast in the spatial frequency domain," *JBO*, vol. 19, p. 056013, 2014.

[73] A. Venkitaraman, S. Chatterjee, and P. Händel, "On Hilbert transform, analytic signal, and modulation analysis for signals over graphs," *Signal Processing*, vol. 156, pp. 106–115, 2019.

[74] L. M. Sanchez-Brea and F. J. Torcal-Milla, "Near-field diffraction of gratings with surface defects," *Appl. Opt.*, vol. 49, pp. 2190–2197, 2010.

[75] X. Zhou, M. Lei, D. Dan, B. Yao, J. Qian, S. Yan, Y. Yang, J. Min, T. Peng, T. Ye, and G. Chen, "Double-exposure optical sectioning structured illumination microscopy based on Hilbert transform reconstruction," *PLOS ONE*, vol. 10, p. e0120892, 2015.

# Supplementary Material for

# *"Fast and Accurate Amplitude Demodulation of Wideband Signals"*

Mantas Gabrielaitis

mantas.gabrielaitis@{ist.ac.at, gmail.com}

## CONTENTS

## OVERVIEW

This document serves as a source of additional information to support the ideas and results introduced in the accompanying paper *"Fast and Accurate Amplitude Demodulation of Wideband Signals."*

Virtually, the material provided in this supplement can be divided into five blocks comprising, respectively, **Sections A – C**, **D – G**, **H – I**, **J – K**, and **L – O**:

- The **A – C** block provides proofs of modulator recovery conditions.
- The **D – G** block is concerned with mathematical aspects of the alternating projection (AP) algorithms of amplitude demodulation.
- The **H – I** block reviews the main types of amplitude-modulated wideband signals found in practice and defines synthetic modulators and carriers used for testing purposes in the present work.
- The **J – K** block presents details on the numerical implementation of the AP and other relevant demodulation methods on a computer, as well as their benchmarking configurations.
- The **L – O** block contains auxiliary simulation results and their discussion.

A summary of each of the sections follows next to ease navigation through this document.

**Section A** establishes several auxiliary results that are exploited in the modulator recovery proofs next.

**Section B** provides proofs for the modulator recovery conditions introduced in Section II-C of the main text.

**Section C** discusses the results and implementation of the numerical experiments performed to extend the modulator recovery conditions.

**Section D** introduces the basic concepts of mathematical analysis necessary for the formulation and study of the properties of AP algorithms.

**Section E** formulates and proves relevant properties of the constraint sets of modulator estimates and defines operators that implement metric projections onto the modulator constraint sets used in this work.

**Section F** provides the convergence proofs of the AP algorithms introduced in the main text.

**Section G** derives a lower bound on the convergence error.

**Section H** reviews the main types of amplitude-modulated wideband signals found in practice.

**Section I** defines the modulators and carriers of synthetic signals used to test amplitude demodulation algorithms in the present work.

**Section J** lists configurations of the execution control parameters employed for the performance analysis of the AP and LDC algorithms of amplitude demodulation.

**Section K** provides information on the implementation and execution of the demodulation algorithms on a computer that we used.

**Section L** overviews the results of demodulation algorithm performance tests in terms of carrier estimates.

**Section M** presents additional simulation results on the dependence of convergence rates of the AP algorithms on the signal length.

**Section N** introduces simulation results of repetitive demodulation of carrier and modulator estimates obtained using the AP-A algorithm.

**Section O** presents the results of additional simulations that demonstrate the suitability and consistency of the AP approach to demodulate wideband speech signals.

## RECOVERY CONDITIONS

### A. AUXILIARY PROOFS

Here, we establish two important properties of the unitary DFT basis vectors that are repetitively used in the proofs of propositions about the modulator recovery conditions in the next section.

**Proposition A.1.** *Consider a subset of DFT basis vectors* $\{\mathbf{f}^{(k)}\}_{k \in \mathcal{I}_n^{\omega^*}}$*. Assume a set of arbitrarily chosen* $n_s$ *sample points encoded by components of a vector* $\mathbf{r} \in \mathbb{N}_+^{n_s}$*, and introduce a linear transform* $\mathbf{L_r}$ *that maps every* $\mathbf{x} \in \mathbb{R}^n$ *to* $[x_{r_1}, x_{r_2}, \ldots, x_{r_{n_s}}]^T \in \mathbb{R}^{n_s}$*. Then, a set* $\{\mathbf{L_r}\mathbf{f}^{(k)}\}_{k \in \mathcal{I}_n^{\omega^*}}$ *is linearly independent if and only if* $n_s \geqslant 2\omega^* - 1$.

*Proof.*

Necessity. If $n_s < 2\omega^* - 1$, the set of vectors $\{\mathbf{L_r}\mathbf{f}^{(k)}\}_{k \in \mathcal{I}_n^{\omega^*}}$ is linearly dependent because the number of linearly independent vectors cannot be higher than the number of components they consists of.

Sufficiency. First, consider the case of $n_s = 2\omega^* - 1$. Then, $\{\mathbf{L_r}\mathbf{f}^{(k)}\}_{k \in \mathcal{I}_n^{\omega^*}}$ is linearly independent if and only if the determinant of a matrix formed by concatenating all vectors from this set in an arbitrary order is not equal to zero (see, e.g., [1, p. 13]). To show that the latter condition is satisfied in our case, define a matrix

$$\mathbf{M} = \mathbf{L_r}[\mathbf{f}^{(n-\omega^*+2)}, \mathbf{f}^{(n-\omega^*+3)}, \ldots, \mathbf{f}^{(n)}, \mathbf{f}^{(1)}, \mathbf{f}^{(2)}, \ldots, \mathbf{f}^{(\omega^*-1)}, \mathbf{f}^{(\omega^*)}]. \tag{24}$$

Taking into account that $\mathbf{L_r}\mathbf{f}^{(k)}$ can be written as $[(z^{r_1})^{k-1}, (z^{r_2})^{k-1}, \ldots, (z^{r_{n_s}})^{k-1}]^T/\sqrt{n}$, with $z = e^{\imath 2\pi/n}$, and that $(z^{r_1})^k = (z^{r_1})^{k \bmod n}$, $\mathbf{M}$ can be expressed as a product of a diagonal matrix and a Vandermonde matrix:

$$\mathbf{M} = \mathrm{diag}[\sqrt{n} \cdot \mathbf{L_r}\mathbf{f}^{(n-\omega^*+2)}] \, [\mathbf{L_r}\mathbf{f}^{(1)}, \mathbf{L_r}\mathbf{f}^{(2)}, \ldots, \mathbf{L_r}\mathbf{f}^{(2\omega^*-1)}]. \tag{25}$$

Thus,

$$\begin{aligned}
\det \mathbf{M} &= \sqrt{n} \cdot \det \mathrm{diag}[\mathbf{L_r}\mathbf{f}^{(n-\omega^*+2)}] \cdot \det[\mathbf{L_r}\mathbf{f}^{(1)}, \mathbf{L_r}\mathbf{f}^{(2)}, \ldots, \mathbf{L_r}\mathbf{f}^{(2\omega^*-1)}] \\
&= \sqrt{n} \cdot \prod_{i=1}^{2\omega^*-1} \left(\mathbf{L_r}\mathbf{f}^{(n-\omega^*+2)}\right)_i \cdot \prod_{i=1}^{2\omega^*-1} \prod_{j=1}^{i-1} \left((\mathbf{L_r}\mathbf{f}^{(2)})_i - (\mathbf{L_r}\mathbf{f}^{(2)})_j\right) \\
&= \frac{1}{\sqrt{n}} \cdot \prod_{i=1}^{2\omega^*-1} \underbrace{e^{\imath 2\pi(n-\omega^*+1)r_i/n}}_{\neq 0} \cdot \prod_{i=2}^{2\omega^*-1} \prod_{j=1}^{i-1} \underbrace{\left(e^{\imath 2\pi r_i/n} - e^{\imath 2\pi r_j/n}\right)}_{\neq 0} \\
&\neq 0,
\end{aligned} \tag{26}$$

which implies that the set $\{\mathbf{L_r}\mathbf{f}^{(k)}\}_{k \in \mathcal{I}_n^{\omega^*}}$ is linearly independent for $n_s = 2\omega^* - 1$. When writing the second equality above, we used the expression of the determinant of a Vandermonde matrix (see, e.g., [1, p. 143]).

It follows from the definition of linear independence that extending each vector in the set by additional components cannot change the set from linearly independent to linearly dependent. Hence, the linear independence of $\{\mathbf{L_r}\mathbf{f}^{(k)}\}_{k \in \mathcal{I}_n^{\omega^*}}$ for $n_s = 2\omega^* - 1$ implies the linear independence of $\{\mathbf{L_r}\mathbf{f}^{(k)}\}_{k \in \mathcal{I}_n^{\omega^*}}$ for $n_s > 2\omega^* - 1$. ∎

*Remark.* The fact that $\{\mathbf{L_r f}^{(k)}\}_{k \in \mathcal{I}_n^{\omega^*}}$ is linearly independent for $n_s \geqslant 2\omega^* - 1$ means that the system of linear equations $\mathbf{L_r x} = \sum_{k \in \mathcal{I}_n^{\omega^*}} \left( a_k \cdot \mathbf{L_r f}^{(k)} \right)$ has a unique solution for all $\mathbf{x} \in \mathcal{S}_{\omega^*}$ if $n_s \geqslant 2\omega^* - 1$. In other words, every $\mathbf{x} \in \mathcal{S}_{\omega^*}$ can be recovered from its known sample points then.

**Proposition A.2.** *Consider some* $\mathbf{x}^{(1)} \in \mathcal{S}_{\omega^*}$ *and* $\mathbf{x}^{(2)} \in \mathcal{S}_{\omega^*}$. *Assume a set of regularly spaced* $n_s \geqslant 2\omega^* - 1$ *sample points encoded by components of a vector* $\mathbf{r} \in \mathbb{N}_+^{n_s}$ *such that* $r_{i+1} - r_i = n/n_s$ *for every* $i \in \mathcal{I}_{n_s}$. *Then,*

$$\|\mathbf{x}^{(1)}\|_2/\|\mathbf{x}^{(2)}\|_2 = \|\mathbf{L_r x}^{(1)}\|_2/\|\mathbf{L_r x}^{(2)}\|_2. \tag{27}$$

*Proof.* By the definition of $\mathcal{S}_{\omega^*}$ (see Section II-A in the main text),

$$\mathbf{x} = \sum_{k \in \mathcal{I}_n^{\omega^*}} \left( a_k \cdot \mathbf{f}^{(k)} \right), \quad \mathbf{x} \in \mathcal{S}_{\omega^*}. \tag{28}$$

Moreover, it follows from the unitary property of the DFT matrix, $\langle \mathbf{f}^{(j)}, \mathbf{f}^{(k)} \rangle = \delta_{j,k}$,[1] that

$$\|\mathbf{x}\|_2^2 = \sum_{k \in \mathcal{I}_n^{\omega^*}} |a_k|^2. \tag{29}$$

Applying $\mathbf{L_r}$ to both sides of (28) yields

$$\mathbf{L_r x} = \sum_{k \in \mathcal{I}_n^{\omega^*}} \left( a_k \cdot \mathbf{L_r f}^{(k)} \right), \quad \mathbf{x} \in \mathcal{S}_{\omega^*}, \tag{30}$$

where, taking into account that $r_{i+1} - r_i = n/n_s$,

$$\begin{aligned} \mathbf{L_r f}^{(k)} &= [e^{\imath 2\pi \cdot r_1 \cdot (k-1)/n}, e^{\imath 2\pi \cdot r_2 \cdot (k-1)/n}, \ldots, e^{\imath 2\pi \cdot r_{n_s} \cdot (k-1)/n}]^{\mathrm{T}}/\sqrt{n} \\ &= [1, e^{\imath 2\pi \cdot 1 \cdot (k-1)/n_s}, \ldots, e^{\imath 2\pi \cdot (n_s-1) \cdot (k-1)/n_s}]^{\mathrm{T}}/\sqrt{n_s} \cdot \left( e^{\imath 2\pi \cdot r_1 \cdot (k-1)/n} \cdot \sqrt{n_s/n} \right). \end{aligned} \tag{31}$$

Note that $\mathbf{L_r f}^{(k)}$ is the $k$-th column of the unitary $n_s \times n_s$ DFT matrix multiplied by a coefficient whose absolute value is equal to $\sqrt{n_s/n}$. Therefore, analogously to (29),

$$\|\mathbf{L_r x}\|_2^2 = (n_s/n) \cdot \sum_{k \in \mathcal{I}_n^{\omega^*}} |a_k|^2, \tag{32}$$

as long as $n_s \geqslant 2\omega^* - 1$.[2] Consequently, $\|\mathbf{L_r x}\|_2 = \sqrt{n_s/n} \cdot \|\mathbf{x}\|_2$ for $\mathbf{x} \in \mathcal{S}_{\omega^*}$, which implies (27). ∎

## B. RECOVERY PROOFS

In this section, we prove the modulator recovery conditions stated in the main text in the form of *Propositions II.1 − II.4*. We repeat the original assertions from the main text for the sake of convenience.

---

[1] Here, and in the sequel, $\delta_{i,j}$ denotes the Kronecker delta.
[2] If $n_s < 2\omega^* - 1$, some of the vectors $\mathbf{L_r f}^{(k)}$ and $\mathbf{L_r f}^{(j)}$ are identical for $k \neq j$, and hence, (32) does not apply.

**Proposition II.1.** *For almost every* $\mathbf{m} \in \mathcal{M}_\omega$, $\hat{\mathbf{m}} = \mathbf{m}$ *only if* $\varpi \geqslant \omega$, *and* $\mathbf{c} \in \mathcal{C}_d$ *with* $n_s \equiv \sum_{i=1}^{n} I_{\{1\}}(|c_i|) \geqslant \varpi + \omega - 1 \implies d \leqslant n - (\varpi + \omega - 2).$[3,4]

*Proof.* We prove the proposition by showing that, almost everywhere in $\mathcal{M}_\omega$, $\mathbf{m} \neq \hat{\mathbf{m}}$ if $\mathbf{c} \notin \mathcal{C}_d$, or $\varpi < \omega$, or $n_s < \varpi + \omega - 1$. For the sake of convenience, we restate the definition of $\hat{\mathbf{m}}$ here:

$$\hat{\mathbf{m}} = \underset{\mathbf{x} \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}{\arg\min} \|\mathbf{x}\|_2. \tag{33}$$

If $\mathbf{c} \notin \mathcal{C}_d$, it means that either $|c_i| < 1$, for every $i \in \mathcal{I}_n$, or there exists at least one $i \in \mathcal{I}_n$ such that $|c_i| > 1$. In the former case, $|s_i|/m_i < 1$ for every $i \in \mathcal{I}_n$. Hence, for $\alpha = \max\{|s_i|/m_i\}_{i \in \mathcal{I}_n}$, $\alpha \cdot \mathbf{m}$ belongs to $\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi$ but has a smaller norm than $\mathbf{m}$, i.e., $\mathbf{m} \neq \hat{\mathbf{m}}$. In the latter case,

$$\left( \underset{\mathbf{x} \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}{\arg\min} \|\mathbf{x}\|_2 \right)_i > m_i \tag{34}$$

for all $i$ corresponding to $|c_i| > 1$ because $|s_i| > m_i$ then. Thus, $\mathbf{m} = \hat{\mathbf{m}}$ does not apply either. Therefore, $\mathbf{m} = \hat{\mathbf{m}}$ holds only if $\mathbf{c} \in \mathcal{C}_d$.

If $\varpi < \omega$, then the subset of modulators for which $\mathbf{m} = \hat{\mathbf{m}}$ is valid has the cardinality of $\mathbb{R}^{2\varpi - 1}$, and hence, has zero volume in $\mathcal{M}_\omega$, whose cardinality is that of $\mathbb{R}^{2\omega - 1}$.

Next, assume that $\mathbf{c} \in \mathcal{C}_d$, and $\varpi \geqslant \omega$, but $n_s < \varpi + \omega - 1$. Let us represent indexes of all sample points corresponding to $|c_i| = 1$ by a vector $\mathbf{r} \in \mathbb{N}_+^{n_s}$. Then, analogously to (30), we have

$$\mathbf{L_r m} = \sum_{k \in \mathcal{I}_n^\varpi} \left( a_k \cdot \mathbf{L_r f}^{(k)} \right). \tag{35}$$

For the sake of convenience, let us redefine

$$\mathbf{f}^{(k)} = \begin{cases} \boldsymbol{\varphi}^{(1)}, & k = 1 \\ (\boldsymbol{\varphi}^{(2k-2)} + \imath \cdot \boldsymbol{\varphi}^{(2k-1)})/\sqrt{2}, & 2 \leqslant k \leqslant \varpi \\ (\boldsymbol{\varphi}^{(2(n-k)+2)} + \imath \cdot \boldsymbol{\varphi}^{(2(n-k)+3)})/\sqrt{2}, & n - \varpi + 2 \leqslant k = n \end{cases}, \tag{36}$$

and

$$a_k = \begin{cases} \alpha_1, & k = 1 \\ (\alpha_{2k-2} + \imath \cdot \alpha_{2k-1})/\sqrt{2}, & 2 \leqslant k \leqslant \varpi \\ (\alpha_{2(n-k)+2} + \imath \cdot \alpha_{2(n-k)+3})/\sqrt{2}, & n - \varpi + 2 \leqslant k = n \end{cases}. \tag{37}$$

Then, (35) turns into

$$\mathbf{L_r m} = \sum_{i=1}^{2\varpi - 1} \left( \alpha_i \cdot \mathbf{L_r \varphi}^{(i)} \right), \tag{38}$$

---

[3]Here $\hat{\mathbf{m}}$ is as defined by (6) in the main text.

[4]As $d \leqslant n - (\varpi + \omega - 2)$ is implied by $\sum_{i=1}^{n} I_{\{1\}}(|c_i|) \geqslant \varpi + \omega - 1$, we do not refer to it explicitly in this proof.

According to *Proposition A.1*, a set of $n_s$ vectors $\{\mathbf{L_r f}^{(1)}, \ldots, \mathbf{L_r f}^{(\lceil (n_s+1)/2 \rceil)}, \mathbf{L_r f}^{(n - \lfloor (n_s - 3)/2 \rfloor)}, \ldots, \mathbf{L_r f}_n\}$ is linearly independent. Hence, by (36), the same applies to $\{\mathbf{L_r} \boldsymbol{\varphi}^{(i)}\}_{i=1}^{n_s}$. Consequently, (38) is an underdetermined system of linear equations defined by a full-rank matrix. We know from linear algebra that a general solution of such system is expressed as a sum of its separate solution $\boldsymbol{\alpha}^{(0)}$ and a solution of

$$\mathbf{0} = \sum_{i=1}^{2\varpi - 1} \left( \alpha_i \cdot \mathbf{L_r} \boldsymbol{\varphi}_i \right), \tag{39}$$

Solutions of (39) form a $(2\varpi - 1 - n_s)$-dimensional subspace of $\mathbb{R}^{2\varpi - 1}$. Thus, we can express the general solution of (38) as

$$\boldsymbol{\alpha} = \boldsymbol{\alpha}^{(0)} + \sum_{i=1}^{2\varpi - 1 - n_s} z_i \boldsymbol{\rho}^{(i)}, \quad \mathbf{z} \in \mathbb{R}^{2\varpi - 1 - n_s}, \tag{40}$$

where $\{\boldsymbol{\rho}^{(i)}\}_{i=1}^{2\varpi - 1 - n_s}$ is an orthonormal basis of the space of solutions of (39). Taking into account (36)$-$(37) as well as the linear independence of $\{\mathbf{f}^{(k)}\}_{k \in \mathcal{I}_n^\varpi}$ and $\{\boldsymbol{\rho}^{(i)}\}_{i=1}^{2\varpi - 1 - n_s}$, (40) together with (30) define a linear injective function that maps from $\mathbb{R}^{2\varpi - 1 - n_s}$ to $\mathcal{S}_\varpi$:

$$f(\mathbf{z}) = \sum_{j=1}^{2\varpi - 1} \left( \boldsymbol{\alpha}^{(0)} + \sum_{i=1}^{2\varpi - 1 - n_s} z_i \boldsymbol{\rho}^{(i)} \right)_j \boldsymbol{\varphi}^{(j)}, \quad \mathbf{z} \in \mathbb{R}^{2\varpi - 1 - n_s}. \tag{41}$$

The image of $f(\mathbf{z})$ is a subset of those elements of $\mathcal{S}_\varpi$ that coincide with the true modulator $\mathbf{m}$ at entries $\mathbf{r} \in \mathbb{N}_+^{n_s}$. The injective nature of this function guarantees the existence of a unique $\mathbf{z_m} \in \mathbb{R}^{2\varpi - 1 - n_s}$ such that $\mathbf{m} = f(\mathbf{z_m})$.

Now, assume an $\mathbf{m} \in \mathcal{M}_\omega^+ = \{\mathbf{x} \in \mathcal{M}_\omega : m_i > 0, i \in \mathcal{I}_n\}$, i.e., any feasible modulator whose all entries are strictly positive. Define an $\epsilon = \min\{m_i - |s_i| : (|c_i| < 1) \wedge (i \in \mathcal{I}_n)\}$. $\epsilon$ exists and is positive because $|c_i| = 1$ only for $n_s < n$ out of $n$ components of $\mathbf{c}$ by the assumption of the proposition, and $m_i - |s_i| = m_i(1 - c_i)$. The linearity of $f(\mathbf{z})$ implies its continuity at every point of its domain. Hence, there exists an $\eta$ such that $\|f(\mathbf{z}) - f(\mathbf{z_m})\|_2 < \epsilon$ whenever $\mathbf{z} \in \mathcal{H}_{\mathbf{z_m}, \eta} = \{\mathbf{z} \in \mathbb{R}^{2\varpi - 1 - n_s} : \|\mathbf{z} - \mathbf{z_m}\|_2 < \eta\}$. On the other hand, $\|f(\mathbf{z}) - f(\mathbf{z_m})\|_2 = \sqrt{\sum_{i=1}^n [(f(\mathbf{z}))_i - m_i]^2} < \epsilon$ implies $|(f(\mathbf{z}))_i - m_i| < \epsilon$, and hence $(f(\mathbf{z}))_i > |s_i|$, for every $i \in \mathcal{I}_n$. Consequently,

$$f(\mathbf{z}) \in (\mathcal{S}_{\geqslant |\mathbf{s}|} \cap \mathcal{S}_\varpi), \quad \mathbf{z} \in \mathcal{H}_{\mathbf{z_m}, \eta}. \tag{42}$$

Furthermore, if $\mathbf{m} = \hat{\mathbf{m}}$, then $\|\mathbf{m}\|_2^2 < \|f(\mathbf{z})\|_2^2$ for every $\mathbf{z} \in (\mathcal{H}_{\mathbf{z_m}, \eta} \backslash \mathbf{z_m})$, i.e., $\mathbf{z_m}$ is a strict local minimum point of

$$\begin{aligned}
\|f(\mathbf{z})\|_2^2 &= \left\| \sum_{j=1}^{2\varpi - 1} \left( \boldsymbol{\alpha}^{(0)} + \sum_{i=1}^{2\varpi - 1 - n_s} z_i \boldsymbol{\rho}^{(i)} \right)_j \boldsymbol{\varphi}^{(j)} \right\|_2^2 \\
&= \left\| \sum_{j=1}^{2\varpi - 1} \left( \boldsymbol{\alpha}^{(0)} + \sum_{i=1}^{2\varpi - 1 - n_s} z_i \boldsymbol{\rho}^{(i)} \right) \right\|_2^2 \\
&= \|\boldsymbol{\alpha}^{(0)}\|_2^2 + \sum_{i=1}^{2\varpi - 1 - n_s} \left( z_i^2 + 2 \cdot z_i \cdot \langle \boldsymbol{\alpha}^{(0)}, \boldsymbol{\rho}^{(i)} \rangle \right).
\end{aligned} \tag{43}$$

$\|f(\mathbf{z})\|_2^2$ is a continuous, differentiable function with a positive-definite Hessian: $\partial^2\|f(\mathbf{z})\|_2^2/\partial z_i\partial z_j = \delta_{i,j}$. Thus, it has a unique strict local minimum point defined by $\partial\|f(\mathbf{z})\|_2^2/\partial z_i = 0$: [5]

$$z_i^\dagger = -\langle \boldsymbol{\alpha}^{(0)}, \boldsymbol{\rho}^{(i)}\rangle, \quad i \in \mathcal{I}_{2\varpi-1-n_s}. \tag{44}$$

Remember that, without loss of generality, $\boldsymbol{\alpha}^{(0)}$ is a particular solution of (38). $\boldsymbol{\alpha}^{(0)}$ corresponding to $\mathbf{m}$, i.e., $\boldsymbol{\alpha}^{(0)} \equiv (\alpha_1^{(0)}, \alpha_2^{(0)}, \ldots, \alpha_{2\omega-1}, 0, \ldots, 0)^{\mathrm{T}}$ with $\mathbf{m} = \sum_{i=1}^{2\omega-1}\alpha_i^{(0)}\boldsymbol{\varphi}^{(i)}$, is exactly such a solution. In this case, as follows from (41), $\mathbf{m} = f(\mathbf{z}^\dagger)$ if and only if $\mathbf{z}^\dagger = \mathbf{0}$.[6] According to (44), that is equivalent to requiring $\boldsymbol{\alpha}^{(0)}$ to be a solution of the following homogeneous system of linear equations:

$$\langle \boldsymbol{\alpha}^{(0)}, \boldsymbol{\rho}^{(i)}\rangle = 0, \quad i \in \mathcal{I}_{2\varpi-1-n_s}. \tag{45}$$

Hence, the subset of $\mathcal{M}_\omega^+$ to which $\mathbf{m} = \hat{\mathbf{m}}$ applies has the same cardinality as $\mathbb{R}^D$, where $D$ is the dimension of the solution space of (45). Taking into account the linear independence of $\{\boldsymbol{\rho}^{(i)}\}_{i=1}^{2\varpi-1-n_s}$, $D$ is equal to the difference between the number of elements of $\boldsymbol{\alpha}^{(0)}$ that are not identically equal to zero and the number of equations that are not trivially satisfied by any feasible $\boldsymbol{\alpha}^{(0)}$. The latter depends on $\mathbf{c}$, specifically, on the positions of sample points with $|c_i| = 1$. To see this, consider two cases:

- A carrier with equidistantly-spaced true sample points: $|c_{i+(j-1)\cdot d}| = 1$ for some $i \in \mathcal{I}_d$ and every $j \in \mathcal{I}_{n_s}$, where $d = n/n_s$. Then, it follows from (31) that some of the elements of the system $\{\mathbf{L_r}f^{(k)}\}_{k\in\mathcal{I}_n^\varpi}$ are identical as long as $n_s < 2\varpi - 1$.[7] Specifically, we have

$$\mathbf{L_r}f^{(k)} = \begin{cases} \mathbf{L_r}f^{(k+n-n_s)}, & \lfloor(n_s+4)/2\rfloor \leqslant k \leqslant \varpi \\ \mathbf{L_r}f^{(k+n_s-n)}, & n - \varpi + 2 \leqslant k \leqslant n - \lfloor(n_s-1)/2\rfloor \end{cases}. \tag{46}$$

  Equivalently,

$$\mathbf{L_r}\boldsymbol{\varphi}_{n_s+\chi_{n_s}-i} = (-1)^{\chi_i} \cdot \mathbf{L_r}\boldsymbol{\varphi}_{n_s-\chi_{n_s}+2\chi_i+i}, \qquad 1 \leqslant i \leqslant 2\varpi - 2 + \chi_{n_s} - n_s \tag{47}$$

  and $\mathbf{L_r}\boldsymbol{\varphi}_{(n_s+1)} = \mathbf{0}$ when $\chi_{n_s} = 0$, where, $\chi_z = z \bmod 2$. Consequently,[8]

$$(\boldsymbol{\rho}^{(i)})_j = \begin{cases} 1/\sqrt{2}, & j = n_s + \chi_{n_s} - i \\ (-1)^{\chi_i+1}/\sqrt{2}, & j = n_s - \chi_{n_s} + 2\chi_i + i, \qquad 1 \leqslant i \leqslant 2\varpi - 2 + \chi_{n_s} - n_s, \\ 0, & \text{otherwise} \end{cases} \tag{48}$$

  and $\boldsymbol{\rho}^{(2\varpi-1-n_s)} = \mathbf{0}$ when $\chi_{n_s} = 0$. Moreover, by our choice of $\boldsymbol{\alpha}^{(0)}$,

$$\alpha_j^{(0)} = 0, \qquad 2\omega \leqslant j \leqslant 2\varpi - 1. \tag{49}$$

---

[5] This strict local minimum point is the only local minimum point of this function.
[6] Note that $\{\boldsymbol{\rho}_i\}_{i=1}^{2\varpi-1-n_s}$ is linearly independent.
[7] Note that $n_s < \varpi + \omega - 1$ and $\varpi \geqslant \omega$ imply $n_s < 2\varpi - 1$.
[8] Note that, as discussed before, $\{\boldsymbol{\varphi}_i\}_{i=1}^{n_s}$ is linearly independent.

(48) and (49) together imply that (45) holds for $i \in \mathcal{I}_{n_s - (2\omega - 1)}$ independent of $\boldsymbol{\alpha}^{(0)}$ corresponding to the chosen $\mathbf{m} \in \mathcal{M}_\omega^+$, and that (45) holds for the remaining $i \in (\mathcal{I}_{2\varpi - 1 - n_s} \setminus \mathcal{I}_{n_s - (2\omega - 1)})$ if and only if

$$\alpha_j^{(0)} = 0, \qquad 2(n_s + 1 - \varpi) \leqslant j \leqslant 2\omega - 1. \tag{50}$$

Hence, (45) applies only to the subset of $\mathcal{M}_\omega^+$ with

$$D = \min\{(2\omega - 1), \ (2\omega - 1) - 2(\varpi + \omega - 1 - n_s)\}$$
$$= \min\{(2\omega - 1), \ 2(n_s + 1 - \varpi) - 1\}. \tag{51}$$

(51) implies that $D < 2\omega - 1$ as long as $n_s < \varpi + \omega - 1$. Then, the subset of all $\mathbf{m} \in \mathcal{M}_\omega^+$ that satisfy $\mathbf{m} = \hat{\mathbf{m}}$ has zero volume in $\mathcal{M}_\omega^+$, which has the cardinality of $\mathbb{R}^{2\omega - 1}$. Moreover, if $n_s < \varpi$, it follows from (50) that (45) has only the trivial solution $\boldsymbol{\alpha}^{(0)} = 0$, which implies $\mathbf{m} = 0$. The latter is infeasible in our case, and hence, $\mathbf{m} = \hat{\mathbf{m}}$ does not apply to any $\mathbf{m} \in \mathcal{M}_\omega^+$. On the other hand, if $n_s \geqslant \varpi + \omega - 1$, then all equations of the (45) system are satisfied independent of the actual feasible $\boldsymbol{\alpha}^{(0)}$. Consequently, $D = 2\omega - 1$.

- An arbitrary $\mathbf{c} \in \mathcal{C}_d$ with unspecified structure. Then, none of the $2\varpi - 1 - n_s$ equations of the system (45) are satisfied independent of the actual $\boldsymbol{\alpha}^{(0)}$.[9] Therefore,

$$D = (2\omega - 1) - (2\varpi - 1 - n_s)$$
$$= n_s - 2(\varpi - \omega). \tag{52}$$

(52) is valid only if $(2\omega - 1) > (2\varpi - 1 - n_s)$. Otherwise, (45) has only the trivial solution $\boldsymbol{\alpha}^{(0)} = \mathbf{0}$, which implies $\mathbf{m} = 0$. The latter is infeasible in our case, and hence, $\mathbf{m} = \hat{\mathbf{m}}$ does not apply to any $\mathbf{m} \in \mathcal{M}_\omega^+$. In any case, $D < (2\omega - 1)$ as long as $n_s < (2\varpi - 1)$, which follows from the assumption of the proposition that $n_s < (\varpi + \omega - 1)$ and $\varpi \geqslant \omega$. In fact, $(2\varpi - 1) - (\varpi + \omega - 1) = (\varpi - \omega)$.

Compared with the general case, a smaller number of necessary sample points with $|c_i| = 1$ is achieved in the first example due to the fact that the subset of vectors $\mathbf{L_r}\boldsymbol{\varphi}^{(k)}$ with $2\omega \leqslant k \leqslant 2\varpi - 1$ is linearly dependent and that $\alpha_k^{(0)}$ are identically equal to zero in that range. Specifically, every $\boldsymbol{\varphi}^{(k)}$ in the range $2\omega \leqslant n_s - 1 + \chi_{n_s}$ is proportional to one of the $\boldsymbol{\varphi}^{(k)}$ in the range $n_s + 2 - \chi_{n_s} \leqslant 2\varpi - 1$. Every such dependence reduces the necessary number of points with $|c_i| = 1$ for modulator recovery by one. Further, it follows from (31), (36), and *Proposition A.1* that the sets $\{\mathbf{L_r}\boldsymbol{\varphi}^{(k)}\}_{k=2\omega}^{n_s - 1 + \chi_{n_s}}$ and $\{\mathbf{L_r}\boldsymbol{\varphi}^{(k)}\}_{k=n_s + 2 - \chi_{n_s}}^{2\varpi - 1}$ are linearly independent. Hence, no further linear dependencies among $\{\mathbf{L_r}\boldsymbol{\varphi}^{(k)}\}_{k=2\omega}^{2\varpi - 1}$ are possible in general. This means that $n_s = \varpi + \omega - 1$ is the absolute minimum of sample points with $|c_i| = 1$ necessary for $\mathbf{m} = \hat{\mathbf{m}}$ to hold. The second example above illustrates that this number is surely higher for some $\mathbf{c} \in \mathcal{C}_d$.

The last three paragraphs demonstrate that the subset of $\mathcal{M}_\omega^+$ to which $\mathbf{m} = \hat{\mathbf{m}}$ applies has zero volume in $\mathcal{M}_\omega^+$ if $n_s < \varpi + \omega - 1$. Taking into account that $(\mathcal{M}_\omega \setminus \mathcal{M}_\omega^+)$ has lower cardinality than $\mathcal{M}_\omega^+$ ($\mathbb{R}^{2\omega - 2}$ vs. $\mathbb{R}^{2\omega - 1}$), we conclude that the set of all $\mathbf{m} \in \mathcal{M}_\omega$ that satisfy $\mathbf{m} = \hat{\mathbf{m}}$ has zero volume in $\mathcal{M}_\omega$ if $n_s < \varpi + \omega - 1$.

---

[9]One possible example of this kind is $\mathbf{c}$ with $|c_i| = 1$ for $i \in \mathcal{I}_{n_s}$, and $|c_i| < 1$ for $i \in (\mathcal{I}_n \setminus \mathcal{I}_{n_s})$.

∎

**Proposition II.2.** *Consider* $\mathbf{m} \in \mathcal{M}_\omega$ *and* $\tilde{\mathbf{c}} \in \mathcal{C}_{\tilde{d}}$ *with* $|\tilde{c}_i| = 1$ *for* $i \in \mathcal{J}_n \subseteq \mathcal{I}_n$, *and* $\tilde{c}_i = 0$ *otherwise. If* $\hat{\mathbf{m}} = \mathbf{m}$ *holds for the* $\mathbf{m}$ *and* $\tilde{\mathbf{c}}$, *then it also holds for every pair made of the same* $\mathbf{m}$ *and any* $\mathbf{c} \in \mathcal{C}_d$ *with* $d \leqslant \tilde{d}$ *and* $|c_i| = 1$ *for* $i \in \mathcal{J}_n$.[10]

*Proof.* Denote $\tilde{\mathbf{s}} = \mathbf{m} \circ \tilde{\mathbf{c}}$ and $\mathbf{s} = \mathbf{m} \circ \mathbf{c}$. Now, note that $|c_i| \geqslant |\tilde{c}_i|$ by the condition of the proposition, and hence, $|s_i| \geqslant |\tilde{s}_i|$. Therefore, $\mathcal{S}_{\geqslant|\mathbf{s}|} \subseteq \mathcal{S}_{\geqslant|\tilde{\mathbf{s}}|}$, which implies $(\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi) \subseteq (\mathcal{S}_{\geqslant|\tilde{\mathbf{s}}|} \cap \mathcal{S}_\varpi)$, and, consequently, $\arg\min_{\mathbf{x}\in\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi} \|\mathbf{x}\|_2 \geqslant \arg\min_{\mathbf{x}\in\mathcal{S}_{\geqslant|\tilde{\mathbf{s}}|}\cap\mathcal{S}_\varpi} \|\mathbf{x}\|_2$. According to the proposition, $\arg\min_{\mathbf{x}\in\mathcal{S}_{\geqslant|\tilde{\mathbf{s}}|}\cap\mathcal{S}_\varpi} \|\mathbf{x}\|_2 = \mathbf{m}$. On the other hand, by construction, $\mathbf{m} \in (\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi)$. Thus, $\arg\min_{\mathbf{x}\in\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi} \|\mathbf{x}\|_2 = \mathbf{m}$. ∎

*Remark.* It can be shown by example that the validity of $\hat{\mathbf{m}} = \mathbf{m}^{(1)}$ for some $\mathbf{m}^{(1)} \in \mathcal{M}_\omega$ and $\tilde{\mathbf{c}} \in \mathcal{C}_d$ with $|\tilde{c}_i| = 1$, $i \in \mathcal{J}_n$, does not necessarily imply the validity of $\hat{\mathbf{m}} = \mathbf{m}^{(2)}$ for another $\mathbf{m}^{(2)} \in \mathcal{M}_\omega$ and the same $\tilde{\mathbf{c}}$.

**Proposition II.3.** *Assume* $\mathbf{m} \in \mathcal{M}_\omega$ *and* $\mathbf{c} \in \mathcal{C}_d$ *with* $\varpi \geqslant \omega$. *If, additionally, there exist* $d \in \mathcal{I}_n$ *and* $i \in \mathcal{I}_d$ *such that* $n_s \equiv (n/d) \in \mathbb{N}_+$, $n_s \geqslant \varpi + \omega - 1$, *and* $|c_{i+(j-1)\cdot d}| = 1$ *for every* $j \in \mathcal{I}_{n_s}$, *then* $\hat{\mathbf{m}} = \mathbf{m}$.[10]

*Proof.* Here, we distinguish between two cases: $(\varpi + \omega - 1) \leqslant n_s < (2\varpi - 1)$ and $n_s \geqslant (2\varpi - 1)$.

If $(\varpi + \omega - 1) \leqslant n_s < (2\varpi - 1)$, then it follows from the proof of *Proposition II.1* that $\mathbf{m}$ has the smallest norm among all elements of the image of $f(\mathbf{z})$ defined by (41), i.e, all elements $\mathbf{x} \in \mathcal{S}_\varpi$ that satisfy $\mathbf{L_r x} = \mathbf{L_r m}$. Next, consider a $\mathbf{y} \in \mathcal{S}_\omega$ such that $(\mathbf{L_r y})_i \geqslant (\mathbf{L_r m})_i$ for every $i \in \mathcal{I}_{n_s}$ and $(\mathbf{L_r y})_i > (\mathbf{L_r m})_i$ for at least one $i \in \mathcal{I}_{n_s}$. Then, by (27), $\|\mathbf{y}\|_2 > \|\mathbf{m}\|_2$. Moreover, using the same argumentation as for $\mathbf{m}$, we see that $\mathbf{y}$ has the smallest norm among all elements $\mathbf{x} \in \mathcal{S}_\varpi$ that satisfy $\mathbf{L_r x} = \mathbf{L_r y}$. Hence, $\mathbf{m}$ has smaller norm than any other element $\mathbf{x} \in \mathcal{S}_\varpi$ that satisfies $\mathbf{L_r x} \geqslant \mathbf{L_r m}$. Moreover, $(\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi) \subset \mathcal{S}_\varpi$. Thus, we conclude that $\hat{\mathbf{m}} = \mathbf{m}$ holds.

If $n_s \geqslant (2\varpi - 1)$, then it follows from (27) of *Proposition A.2* that

$$\arg\min_{\mathbf{x}\in\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi} \|\mathbf{x}\|_2 = \arg\min_{\mathbf{x}\in\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi} \|\mathbf{L_r x}\|_2. \tag{53}$$

Further, the constraint set $\mathcal{S}_{\geqslant|\mathbf{s}|}$ implies that

$$\|\mathbf{L_r x}\|_2 \geqslant \|\mathbf{L_r}|\mathbf{s}|\|_2, \quad \mathbf{x} \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi. \tag{54}$$

On the other hand, $|s_{r_i}| = |m_{r_i} \cdot c_{r_i}| = m_{r_i} \cdot |c_{r_i}| = m_{r_i}$ for $i \in \mathcal{I}_{n_s}$, i.e., $\mathbf{L_r}|\mathbf{s}| = [m_{r_1}, m_{r_2}, \ldots, m_{r_{n_s}}]^{\mathrm{T}}$. According to *Proposition A.1*, $\{\mathbf{L_r f}^{(k)}\}_{k\in\mathcal{I}_n^\varpi}$ is linearly independent if $n_s \geqslant 2\varpi - 1$. Therefore, (30)[11] has a unique solution, which, by (28), means that, among all $\mathcal{S}_\varpi$, there is a unique $\mathbf{x} = \mathbf{m}$ that satisfies $\mathbf{L_r x} = [m_{r_1}, m_{r_2}, \ldots, m_{r_{n_s}}]^{\mathrm{T}} = \mathbf{L_r}|\mathbf{s}|$. Hence, by (54), $\arg\min_{\mathbf{x}\in\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi} \|\mathbf{L_r x}\|_2 = \mathbf{m}$, and, by (53), $\arg\min_{\mathbf{x}\in\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi} \|\mathbf{x}\|_2 = \mathbf{m}$. ∎

---

[10]Here $\hat{\mathbf{m}}$ is as defined by (6) in the main text.

[11]Note that $\omega^*$ in (28)–(30) stands for $\varpi$ here.

**Proposition II.4.** *Consider* $\mathbf{m} \in \mathcal{M}_\omega$ *and* $\mathbf{c} \in \mathcal{S}_{|..|\leqslant 1}$. *Take* $n_s \geqslant 2\varpi - 1$ *sample points of* $\mathbf{s} = \mathbf{m} \circ \mathbf{c}$ *whose indexes are defined as entries of any chosen* $\mathbf{r} \in \mathbb{N}_+^{n_s}$ *with* $r_{i+1} - r_i = n/n_s$ *for every* $i \in \mathcal{I}_{n_s}$. *Then,*[12]

$$\|\mathbf{m} - \hat{\mathbf{m}}\|_2 / \|\mathbf{m}\|_2 \leqslant \sqrt{1 - \sum_{i=1}^{n_s} s_{r_i}^2 / \sum_{i=1}^{n_s} m_{r_i}^2}. \tag{55}$$

*Proof.* For the sake of convenience, we will exploit the linear transformation $\mathbf{L_r}$, already introduced in *Proposition A.1*, that maps every $\mathbf{x} \in \mathbb{R}^n$ to $[x_{r_1}, x_{r_2}, \ldots, x_{r_{n_s}}]^{\mathrm{T}}$. Then, (55) can be rewritten as

$$\|\mathbf{m} - \hat{\mathbf{m}}\|_2 / \|\mathbf{m}\|_2 \leqslant \sqrt{1 - \|\mathbf{L_r s}\|_2^2 / \|\mathbf{L_r m}\|_2^2}. \tag{56}$$

Note that

$$\|\mathbf{m}\|_2^2 - \|\hat{\mathbf{m}}\|_2^2 - \|\mathbf{m} - \hat{\mathbf{m}}\|_2^2 = 2 \cdot \|\hat{\mathbf{m}}\|_2 \cdot (\|\mathbf{m}\|_2 - \|\hat{\mathbf{m}}\|_2). \tag{57}$$

Next, we have by construction that $\mathbf{m} \in (\mathcal{S}_{\geqslant |\mathbf{s}|} \cap \mathcal{S}_\varpi)$. Hence,

$$\|\mathbf{m}\|_2^2 \geqslant \|\hat{\mathbf{m}}\|_2^2, \tag{58}$$

which, together with (57) implies $\|\mathbf{m}\|_2^2 - \|\hat{\mathbf{m}}\|_2^2 - \|\mathbf{m} - \hat{\mathbf{m}}\|_2^2 \geqslant 0$, i.e.,

$$\|\mathbf{m} - \hat{\mathbf{m}}\|_2 / \|\mathbf{m}\|_2 \leqslant \sqrt{1 - \|\hat{\mathbf{m}}\|_2^2 / \|\mathbf{m}\|_2^2}. \tag{59}$$

On the other hand, by *Proposition A.2* , $\|\hat{\mathbf{m}}\|_2^2 / \|\mathbf{m}\|_2^2 = \|\mathbf{L_r \hat{m}}\|_2^2 / \|\mathbf{L_r m}\|_2^2$ if $n_s \geqslant 2\varpi - 1$ and $r_{i+1} - r_i = n/n_s$ for every $i \in \mathcal{I}_{n_s}$. Thus,

$$\|\mathbf{m} - \hat{\mathbf{m}}\|_2 / \|\mathbf{m}\|_2 \leqslant \sqrt{1 - \|\mathbf{L_r \hat{m}}\|_2^2 / \|\mathbf{L_r m}\|_2^2}. \tag{60}$$

Finally, $\hat{m}_i \geqslant |s_i|$ for every $i \in \mathcal{I}_n$ because $\mathbf{m} \in \mathcal{S}_{\geqslant |\mathbf{s}|}$, which means that

$$\|\mathbf{L_r \hat{m}}\|_2^2 \geqslant \|\mathbf{L_r s}\|_2^2. \tag{61}$$

Combining (60) with (61) leads to (56).                                                                          ∎

*Remark.* Note that $\|\mathbf{m}\|_2^2 = \|\hat{\mathbf{m}}\|_2^2$ in (58) if and only if $\mathbf{m} = \hat{\mathbf{m}}$ because $\mathcal{S}_{\geqslant |\mathbf{s}|} \cap \mathcal{S}_\varpi$ is convex and $\|\ldots\|_2^2$ is strictly convex. Therefore, the equality in (55) holds if and only if $\mathbf{m} = \hat{\mathbf{m}}$, i.e., the modulator recovery is exact.

## C. FURTHER ANALYSIS: NUMERICAL EXPERIMENTS

Here, we present the results of numerical experiments used to extend the modulator recovery conditions to carriers with nonuniformly placed sample points $|c_i| = 1$ in terms of the parameters $n$, $\omega$, $\varpi$, and $d$.

---

[12]Here $\hat{\mathbf{m}}$ is as defined by (6) in the main text.

*Setup*

The numerical experiments under consideration consist of the following steps.

1. $10^3$ pairs of $\mathbf{m}$ and $\mathbf{c}$ are generated by randomly sampling from $\mathcal{M}_\omega$ and $\mathcal{C}_d$ for every feasible combination of $\omega$ and $d$ consistent with a chosen signal length $n$.

2. For every pair of $\mathbf{m}$ and $\mathbf{c}$ generated, $\mathbf{m}$ is inferred from the $\mathbf{s} = \mathbf{m} \circ \mathbf{c}$ via $\hat{\mathbf{m}}$ defined by (6). The latter is evaluated by using the AP-P algorithm, introduced in Section III-C of the main text, with $\epsilon_{tol} = 10^{-14}$ and unlimited $N_{iter}$.

3. For every combination of the parameters $\omega$ and $d$, two estimates related to the recovery error are evaluated: 1) the average empirical error $\langle E_m \rangle$ and 2) the fraction of cases with vanishing error $P(E_m < \varepsilon)$, where $\varepsilon$ is a positive number arbitrarily close to zero. $P(E_m < \varepsilon)$ can be seen as the demodulation success rate for a given error threshold $\varepsilon$.

A crucial aspect of the above experiments to producing informative data for our purposes is the way $\mathcal{M}_\omega$ and $\mathcal{C}_d$ are sampled. For both of these sets, we exploited uniform sampling but with some additional constraints, as explained next.

- The cutoff frequency $\omega$, defining the modulator set, and the cutoff frequency $\varpi$, defining the estimator $\hat{\mathbf{m}}$, were fixed to be equal. This choice allowed us to considerably reduce the extent of relevant parameter combinations to be checked without loss of generality. Indeed, $\varpi \geqslant \omega$ is a necessary condition for a full recovery independent of $\mathbf{c} \in \mathcal{C}_d$ and $\mathcal{M}_\omega \subset \mathcal{M}_\varpi$ if $\varpi > \omega$. Hence, all recovery conditions applicable in the case of $\varpi = \omega$ hold for $\varpi > \omega$ as well.

- Only the subset of pure spike-train carriers consisting of $c_i \in \{0, 1\}$ sample points was considered among all possible $\mathbf{c} \in \mathcal{C}_d$. According to *Proposition II.2*, that is sufficient for identifying full recovery conditions without loss of generality.

- Different elements of the pure spike-train subset of $\mathcal{C}_d$ may substantially differ in the number $n_s$ of sample points with $|c_i| = 1$. In particular, $\lceil n/d \rceil \leqslant n_s \leqslant n - d + 1$. We considered modulator reconstruction by uniformly sampling either from the sparsest ($n_s = \lceil n/d \rceil$) or the densest ($n_s = n - d + 1$) subset of spike-train carriers. In view of *Proposition II.2*, pure spike-train carriers with $n_s = \lceil n/d \rceil$ have the tightest, and hence the most general, constraints for exact modulator recovery in terms of the parameters $\varpi$ and $d$.

The algorithms for sampling from the modulator and carrier sets specified above are presented in Section I.

*Results*

Fig. 10 $A_1$ displays color-plots of the fraction $P$ of the modulator recovery cases with $E_m < 10^{-12}$ over the $(d, \varpi)$ plane for $n = 32$ and $n_s = \lceil n/d \rceil$. In agreement with the necessary recovery condition discussed in Section II-C, $P(E_m < 10^{-12})$ is equal to 0 for all points with $\lceil n/d \rceil < 2\varpi - 1$. More remarkably, for $\lceil n/d \rceil \geqslant 2\varpi - 1$, $P(E_m < 10^{-12})$ equals 1 except some boundary points, where $P$ varies between 0.65 and 1.
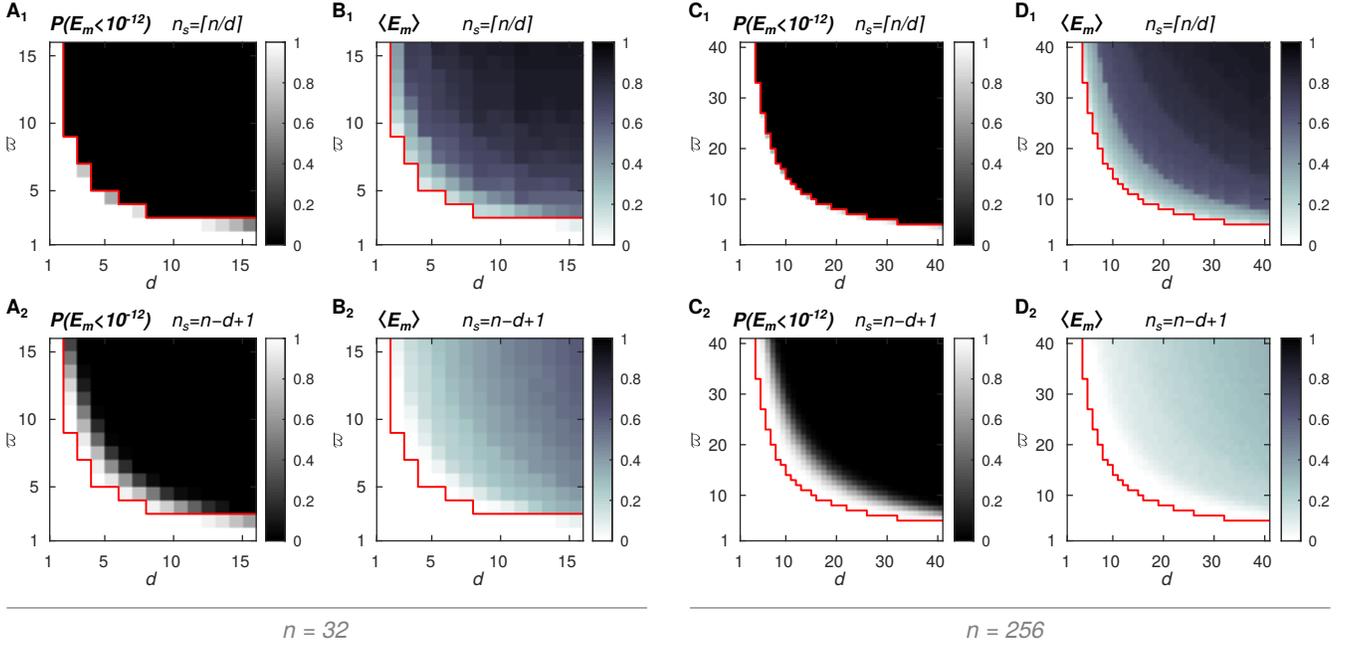
Fig. 10. Success rates and errors of modulator recovery for signals based on pure $c_i \in \{0, 1\}$ spike-train carriers. $\mathbf{A_1}$–$\mathbf{A_2}$: color plots of the fraction ($P$) of modulator recovery cases with the recovery error $E_m$ lower than $10^{-12}$ for different combinations of $d$ and $\varpi$, and $n = 32$; red lines plot the relation $\lceil n/d \rceil = 2\varpi - 1$. $A_1$ displays the results for carriers with $n_s = \lceil n/d \rceil$ spikes, while $A_2$ corresponds to carriers with $n_s = n - d + 1$. $\mathbf{B_1}$–$\mathbf{B_2}$: the same as $A_1$–$A_2$, but with the average recovery error instead of the success rate over all modulator and carrier pairs shown for each combination of $d$ and $\varpi$. $\mathbf{C_1}$–$\mathbf{C_2}$: the same as $A_1$–$A_2$ except that $n = 256$. $\mathbf{D_1}$–$\mathbf{D_2}$: the same as $B_1$–$B_2$ except that $n = 256$.

However, even in the latter cases, the error is small, as follows from Fig. 10 $B_1$, which plots the average $\langle E_m \rangle$ over the $(d, \omega)$ plane. The maximum likelihood (ML) estimate of $P(E_m < 10^{-12})$ for all tested modulator-carrier pairs that adhere to the necessary recovery condition $\lceil n/d \rceil \geqslant 2\varpi - 1$ is 0.971; its 99% confidence interval is $(0.969, 0.972)$.

Increasing the number of carrier points with $|c_i| = 1$ does not change the landscape of the recovery success rate considerably. Indeed, pushing $n_s$ to the maximum $n - d + 1$ increases the $P$ values only at points in the immediate vicinity of the $\lceil n/d \rceil = 2\varpi - 1$ boundary (see Fig. 10 $A_2$). Nevertheless, it has to be noted that the recovery errors are decreased by increasing $n_s$ on average (see Fig. 10 $B_2$). The ML estimate of $P(E_m < 10^{-12})$ for all tested modulator-carrier pairs that adhere to the necessary recovery condition $\lceil n/d \rceil \geqslant 2\varpi - 1$ is 0.987 in this case; its 99% confidence interval is $(0.986, 0.988)$.

We found an analogous picture while considering signals with different lengths $n$. One such example ($n = 256$) is considered in Fig. 10 $C_1$–$D_2$. The chances of the modulator recovery with vanishing error, i.e., $E_m < 10^{-12}$, are higher in this case. Specifically, the ML estimate of $P(E_m < 10^{-12})$ is 0.9933, with the 99% confidence interval $(0.9930, 0.9936)$ when $n_s = \lceil n/d \rceil$. That can be explained by the smaller contribution of the boundary

points of the relation $\lceil n/d \rceil = 2\varpi - 1$ in the $(d, \varpi)$ plane to the total count. It is important to note that, in all cases discussed here, essentially the same results are obtained even if the error threshold $\varepsilon$ is increased to $10^{-3}$. This rejects any possibility of numerical inaccuracies affecting our conclusion.

## AP ALGORITHMS

### D. MATHEMATICAL PRELIMINARIES

In this section, we introduce some basic concepts of mathematical analysis necessary for the formulation and assessment of the AP algorithms that we used to calculate modulator estimators defined by (6) and (8) in the main text.

*Convex, interior, closed, and bounded sets*

We start with definitions of a few basic attributes of sets in Euclidean spaces.

**Definition D.1.** A set $\mathcal{S} \subseteq \mathbb{R}^n$ is said to be *convex* if $\theta \cdot \mathbf{x} + (1 - \theta) \cdot \mathbf{y} \in \mathcal{S}$ for all $\mathbf{x}, \mathbf{y} \in \mathcal{S}$ and $\theta \in [0, 1]$.

**Definition D.2.** An element $\mathbf{x} \in \mathcal{S} \subseteq \mathbb{R}^n$ is said to be *an interior point* of $\mathcal{S}$ if there exists an $\epsilon > 0$ such that $\{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{y}\|_2 < \epsilon\} \subset \mathcal{S}$.

**Definition D.3.** A set that consists of all interior points of $\mathcal{S} \subseteq \mathbb{R}^n$ is called *the interior* of $\mathcal{S}$. We denote it by $\mathcal{S}^\circ$.

**Definition D.4.** An element $\mathbf{y} \in \mathbb{R}^n$ is said to be *a contact point* of $\mathcal{S} \subseteq \mathbb{R}^n$ if, for any $\epsilon > 0$, there exists an $\mathbf{x} \in \mathcal{S}$ such that $\|\mathbf{x} - \mathbf{y}\|_2 < \epsilon$.

**Definition D.5.** A set $\mathcal{S} \subseteq \mathbb{R}^n$ is said to be *closed* if it is equal to the set of all its contact points.

Convexity and closedness of sets are preserved under intersection.

**Proposition D.6.** *The intersection $\mathcal{S}_1 \cap \mathcal{S}_2$ of two closed and convex sets $\mathcal{S}_1$ and $\mathcal{S}_2$ is closed and convex.*

Another important characteristic of sets is their boundedness.

**Definition D.7.** A set $\mathcal{S} \subset \mathbb{R}^n$ is said to be *bounded* if there exists a $b \in \mathbb{R}$ such that $\|\mathbf{x} - \mathbf{y}\|_2 \leqslant b$ for all $\mathbf{x}, \mathbf{y} \in \mathcal{S}$.

**Definition D.8.** A set $\mathcal{S} \subset \mathbb{R}^1 = \mathbb{R}$ is said to be *bounded from above* if there exists a $u \in \mathbb{R}$ such that $\mathbf{x} \leqslant u$ for all $\mathbf{x} \in \mathcal{S}$. $u$ is called an *upper bound* of $\mathcal{S}$.

*Remark.* $\underline{u} \in \mathbb{R}$ is said to be the *least upper bound* of $\mathcal{S} \subset \mathbb{R}$ if $\mathbf{x} \leqslant \underline{u}$ for all $\mathbf{x} \in \mathcal{S}$, and there exists a $\mathbf{y} \in \mathcal{S}$ for every $\epsilon > 0$ such that $\mathbf{y} > \underline{u} - \epsilon$. The least upper bound exists for any $\mathcal{S} \subset \mathbb{R}$ bounded from above due to the continuity of the real numbers.

**Definition D.9.** A set $\mathcal{S} \subset \mathbb{R}^1 = \mathbb{R}$ is said to be *bounded from below* if there exists an $l \in \mathbb{R}$ such that $l \leqslant \mathbf{x}$ for all $\mathbf{x} \in \mathcal{S}$. $l$ is called a *lower bound* of $\mathcal{S}$.

*Remark.* $\bar{l} \in \mathbb{R}$ is said to be the *greatest lower bound* of $\mathcal{S} \subset \mathbb{R}$ if $\bar{l} \leqslant \mathbf{x}$ for all $\mathbf{x} \in \mathcal{S}$, and there exists a $\mathbf{y} \in \mathcal{S}$ for every $\epsilon > 0$ such that $\mathbf{y} < \bar{l} + \epsilon$. Analogously to the least upper bound, any $\mathcal{S} \subset \mathbb{R}$ bounded from below has the greatest lower bound.

### *Convergence of sequences*

The following fundamental properties of infinite sequences of points in bounded subsets of Euclidean spaces play a critical role in the proofs of the convergence of the AP algorithms.

**Proposition D.10** (Monotone Convergence Theorem)**.** *Any monotonically decreasing sequence of real numbers*

$$x^{(0)} \geqslant x^{(1)} \geqslant \ldots \geqslant x^{(i)} \geqslant \ldots \tag{62}$$

*that is bounded from below converges to its greatest lower bound* $\bar{l}$, *i.e., for every* $\epsilon > 0$, *there exists an* $N(\epsilon)$ *such that* $|x^{(i)} - \bar{l}| < \epsilon$ *whenever* $i > N(\epsilon)$.

*Remark.* Analogously, any monotonically increasing sequence that is bounded from above converges to its least upper bound.

**Definition D.11.** Consider a sequence $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(i)}, \ldots$ in $\mathbb{R}^n$, i.e., $\mathbf{x}^{(i)} \in \mathbb{R}^n$ for every $i \geqslant 0$. Another sequence $\mathbf{x}^{(k_0)}, \mathbf{x}^{(k_1)}, \ldots, \mathbf{x}^{(k_i)}, \ldots$ in $\mathbb{R}^n$ generated by removing some of the elements of the original sequence is called a *subsequence* of the latter. Note that $k_i > k_j$ for all $i > j \geqslant 0$, and $k_i \geqslant i$ for every $i \geqslant 0$ here.

**Proposition D.12** (Bolzano-Weierstrass Theorem)**.** *Any bounded infinite sequence* $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(i)}, \ldots$ *in* $\mathbb{R}^n$ *has an infinite subsequence* $\mathbf{x}^{(k_0)}, \mathbf{x}^{(k_1)}, \ldots, \mathbf{x}^{(k_i)}, \ldots$ *that converges to a particular* $\mathbf{x}^\dagger \in \mathbb{R}^n$, *i.e., for any* $\epsilon > 0$, *there exists an* $N(\epsilon)$ *such that* $\|\mathbf{x}^{(k_i)} - \mathbf{x}^\dagger\|_2 < \epsilon$ *whenever* $i > N(\epsilon)$.

### *Metric projections*

The central operation around which AP algorithms are built is that of a metric projection.

**Definition D.13.** An element $\mathbf{x_z}$ of a closed subset $\mathcal{S}$ of $\mathbb{R}^n$ is said to be *a metric projection* of $\mathbf{z} \in \mathbb{R}^n$ onto $\mathcal{S}$ if $\|\mathbf{x_z} - \mathbf{z}\|_2 \leqslant \|\mathbf{x} - \mathbf{z}\|_2$ for all $\mathbf{x} \in \mathcal{S}$. We denote a transformation that assigns an $\mathbf{x_z} \in \mathcal{S}$ to every $\mathbf{z} \in \mathbb{R}^n$ by $\mathbf{P}_\mathcal{S} : \mathbb{R}^n \to \mathcal{S}$.

*Remark.* $\mathbf{P}_\mathcal{S}$ generalizes the linear projection operator that assigns an element of a linear space to one of its subspaces (see, e.g., [2, p. 223]). For the sake of brevity, we skip the qualifier "metric" and refer to $\mathbf{P}_\mathcal{S}$ as "a projection" in the sequel.

**Proposition D.14** (see Theoreom 5.11 in [3])**.** *A projection of any element of* $\mathbb{R}^n$ *onto its closed convex subset* $\mathcal{S}$ *exists and is unique.*

*Inequalities*

Two important inequalities that we use extensively in the convergence proofs of AP algorithms apply to Euclidean spaces.

**Proposition D.15** (Triangle Inequality)**.**

$$\|\mathbf{x} + \mathbf{y}\|_2 \leqslant \|\mathbf{x}\|_2 + \|\mathbf{y}\|_2 \qquad \forall \mathbf{x} \in \mathbb{R}^n, \forall \mathbf{y} \in \mathbb{R}^n. \tag{63}$$

*Remark.* Another inequality relevant to us follows from (63). In particular, let us consider some $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{R}^n$. If we define $\mathbf{x} = \mathbf{a} - \mathbf{b}$ and $\mathbf{y} = \mathbf{b} - \mathbf{c}$, then (63) implies $\|\mathbf{a} - \mathbf{c}\|_2 \leqslant \|\mathbf{a} - \mathbf{b}\|_2 + \|\mathbf{b} - \mathbf{c}\|_2$, i.e., $\|\mathbf{a} - \mathbf{b}\|_2 \geqslant \|\mathbf{c} - \mathbf{a}\|_2 - \|\mathbf{c} - \mathbf{b}\|_2$. On the other hand, setting $\mathbf{x} = \mathbf{c} - \mathbf{a}$ and $\mathbf{y} = \mathbf{a} - \mathbf{b}$, we obtain $\|\mathbf{c} - \mathbf{b}\|_2 \leqslant \|\mathbf{c} - \mathbf{a}\|_2 + \|\mathbf{a} - \mathbf{b}\|_2$, i.e., $\|\mathbf{a} - \mathbf{b}\|_2 \geqslant -(\|\mathbf{c} - \mathbf{a}\|_2 - \|\mathbf{c} - \mathbf{b}\|_2)$. Thus,

$$\|\mathbf{a} - \mathbf{b}\|_2 \geqslant \big| \|\mathbf{c} - \mathbf{a}\|_2 - \|\mathbf{c} - \mathbf{b}\|_2 \big| \qquad \forall \mathbf{a} \in \mathbb{R}^n, \forall \mathbf{b} \in \mathbb{R}^n, \forall \mathbf{c} \in \mathbb{R}^n. \tag{64}$$

**Proposition D.16** (Containing-Half-Space Inequality, see Theoreom 5.13 in [3])**.** *If $\mathcal{S} \subset \mathbb{R}^n$ is closed and convex, then*

$$\langle \mathbf{x} - \mathbf{P}_{\mathcal{S}}[\mathbf{x}], \mathbf{P}_{\mathcal{S}}[\mathbf{x}] - \mathbf{y} \rangle \geqslant 0 \qquad \forall \mathbf{x} \in \mathbb{R}^n, \forall \mathbf{y} \in \mathcal{S}. \tag{65}$$

*Remark.* $\mathcal{S}$ belongs to a half-space $\mathcal{H}_{\mathbf{x}} = \{\mathbf{z} \in \mathbb{R}^n : \langle \mathbf{x} - \mathbf{P}_{\mathcal{S}}[\mathbf{x}], \mathbf{P}_{\mathcal{S}}[\mathbf{x}] - \mathbf{z} \rangle \geqslant 0\}$, which is defined for every $\mathbf{x} \in (\mathbb{R}^n \backslash \mathcal{S})$. $\mathbf{P}_{\mathcal{S}}[\mathbf{x}]$ is a boundary point of the $\mathcal{H}_{\mathbf{x}}$.

## E. PROPERTIES OF THE CONSTRAINT SETS AND ASSOCIATED METRIC PROJECTIONS

In this section, we establish the convexity, closedness, and other relevant properties of the constraint sets $\mathcal{S}_{\geqslant |\mathbf{s}|}$ and $\mathcal{S}_{\varpi}$ that lay the basis for the formulation of the AP demodulation algorithms and determine their convergence properties. We then define the concrete metric projection operators of points in $\mathbb{R}^n$ onto $\mathcal{S}_{\geqslant |\mathbf{s}|}$ and $\mathcal{S}_{\varpi}$, which are the main building blocks of the AP demodulation algorithms defined in Section III of the main text.

*Properties of the constraint sets*

The constraint sets $\mathcal{S}_{\geqslant |\mathbf{s}|}$ and $\mathcal{S}_{\varpi}$ that define the AP approach to demodulation introduced in Section II of the main text have the following properties.

**Proposition E.1.** *The set $\mathcal{S}_{\geqslant |\mathbf{s}|}$ is convex and closed. Its interior is $\mathcal{S}_{\geqslant |\mathbf{s}|}^{\circ} = \{\mathbf{x} \in \mathbb{R}^n : x_i > |s_i|, \ i \in \mathcal{I}_n\}$.*

*Proof.* The range of values of each component $x_i$ of $\mathbf{x} \in \mathcal{S}_{\geqslant |\mathbf{s}|}$ and $y_i$ of $\mathbf{y} \in \mathcal{S}_{\geqslant |\mathbf{s}|}$ is $[|s_i|, +\infty) \subset \mathbb{R}$. Obviously,

$$\theta \cdot x_i + (1 - \theta) \cdot y_i \geqslant \min[x_i, y_i] \geqslant |s_i|$$

and

$$\theta \cdot x_i + (1 - \theta) \cdot y_i \leqslant \max[x_i, y_i] < +\infty,$$

where $\theta \in [0, 1]$. Therefore,

$$\theta \cdot x_i + (1 - \theta) \cdot y_i \in [|s_i|, +\infty) \qquad \forall \theta \in [0, 1], \forall i \in \mathcal{I}_n.$$

This, in turn, implies $\theta \cdot \mathbf{x} + (1 - \theta) \cdot \mathbf{y} \in \mathcal{S}_{\geqslant|\mathbf{s}|}$ because the range of values of each component of $\mathbf{x}$ and $\mathbf{y}$ is independent of the actual values of the remaining components. Hence, by *Definition D.1*, $\mathcal{S}_{\geqslant|\mathbf{s}|}$ is convex.

To show that $\mathcal{S}_{\geqslant|\mathbf{s}|}$ is closed, we first note that its complement $\overline{\mathcal{S}}_{\geqslant|\mathbf{s}|} = \mathbb{R}^n \backslash \mathcal{S}_{\geqslant|\mathbf{s}|}$ is equal to $\{\mathbf{x} \in \mathbb{R}^n : x_i < |s_i|, \ i \in \mathcal{I}_n\}$. For any $\mathbf{x} \in \overline{\mathcal{S}}_{\geqslant|\mathbf{s}|}$ and $\epsilon \leqslant \min[|\mathbf{s}| - \mathbf{x}]$,[13] there exists no $\mathbf{y} \in \mathcal{S}_{\geqslant|\mathbf{s}|}$ such that $\|\mathbf{x} - \mathbf{y}\|_2 < \epsilon$. Thus, none of the elements of $\overline{\mathcal{S}}_{\geqslant|\mathbf{s}|}$ are contact points of $\mathcal{S}_{\geqslant|\mathbf{s}|}$. On the other hand, trivially, all points of $\mathcal{S}_{\geqslant|\mathbf{s}|}$ are its contact points. Hence, $\mathcal{S}_{\geqslant|\mathbf{s}|}$ coincides with the set of its contact points, i.e., it is a closed set.

To determine the interior of $\mathcal{S}_{\geqslant|\mathbf{s}|}$, let us denote

$$\mathcal{A}_0 = \{\mathbf{x} \in \mathbb{R}^n : x_i > |s_i|, \ i \in \mathcal{I}_n\},$$
$$\mathcal{A}_i = \{\mathbf{x} \in \mathbb{R}^n : x_i = |s_i|, x_j \geqslant |s_j|, \ j \in (\mathcal{I}_n \backslash \{i\})\} \qquad \forall i \in \mathcal{I}_n.$$

By construction, $\mathcal{S}_{\geqslant|\mathbf{s}|} = \mathcal{A}_0 \cup (\cup_{i=1}^n \mathcal{A}_i)$. Next, we note that

$$\{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{y}\|_2 < \epsilon\} \subset \mathcal{S}_{\geqslant|\mathbf{s}|} \qquad \forall \epsilon \leqslant \min[\mathbf{x} - |\mathbf{s}|], \ \forall \mathbf{x} \in \mathcal{A}_0.$$

Thus, by *Definition D.2*, all points of $\mathcal{A}_0$ are interior points of $\mathcal{S}_{\geqslant|\mathbf{s}|}$. On the other hand, for all $\mathbf{x} \in \mathcal{A}_i$ and $\epsilon > 0$, there exists a $\mathbf{y} \in \mathbb{R}^n$ such that $\|\mathbf{x} - \mathbf{y}\|_2 < \epsilon$ and $\mathbf{y} \notin \mathcal{S}_{\geqslant|\mathbf{s}|}$. In particular, this is satisfied by $\mathbf{y}$ such that $y_i = x_i - \bar{\epsilon}$ with $0 < \bar{\epsilon} < \epsilon$ and $y_j = x_j$ for all $j \in (\mathcal{I}_n \backslash \{i\})$. Therefore, none of $\mathbf{x} \in \mathcal{A}_i$ with $i \in \mathcal{I}_n$ are interior points of $\mathcal{S}_{\geqslant|\mathbf{s}|}$. Altogether, this allows us to conclude that the interior of $\mathcal{S}_{\geqslant|\mathbf{s}|}$ is equal to $\mathcal{A}_0$. That $\mathcal{A}_0$ is nonempty follows straightforwardly from the fact that there exists an $x_i > s_i$ for any $s_i \in \mathbb{R}$. ∎

**Proposition E.2.** *The set $\mathcal{S}_\varpi$ is convex, closed, and void of interior points. In particular, $\mathcal{S}_\varpi$ is a linear subspace of $\mathbb{R}^n$.*

*Proof.* It follows directly from the definition of $\mathbb{R}^n$ that

$$\theta \cdot \mathbf{x} + (1 - \theta) \cdot \mathbf{y} \in \mathbb{R}^n \qquad \forall \mathbf{x} \in \mathbb{R}^n, \forall \mathbf{y} \in \mathbb{R}^n, \forall \theta \in [0, 1], \tag{66}$$

and thus,[14]

$$\theta \cdot \mathbf{x}_\varpi + (1 - \theta) \cdot \mathbf{y}_\varpi \in \mathbb{R}^n \qquad \forall \mathbf{x}_\omega \in \mathcal{S}_\varpi, \forall \mathbf{y}_\varpi \in \mathcal{S}_\varpi, \forall \theta \in [0, 1]. \tag{67}$$

The definition of $\mathcal{S}_\varpi$ implies that $(\mathbf{F}\mathbf{x}_\varpi)_i = 0$ and $(\mathbf{F}\mathbf{y}_\varpi)_i = 0$ for all $\mathbf{x}_\varpi \in \mathcal{S}_\varpi$, $\mathbf{y}_\varpi \in \mathcal{S}_\varpi$, and $i \in (\mathcal{I}_n \backslash \mathcal{I}_n^\varpi)$, so that

$$\theta \cdot (\mathbf{F}\mathbf{x}_\varpi)_i + (1 - \theta) \cdot (\mathbf{F}\mathbf{y}_\varpi)_i = (\mathbf{F}(\theta \cdot \mathbf{x}_\varpi + (1 - \theta) \cdot \mathbf{y}_\varpi))_i = 0 \qquad \forall i \in (\mathcal{I}_n \backslash \mathcal{I}_n^\varpi) \tag{68}$$

---

[13]Note that $\min[|\mathbf{s}| - \mathbf{x}] > 0$.

[14]Note that $\mathcal{S}_\varpi \subset \mathbb{R}^n$.

because of the linearity of the Fourier transform. Combining (67) and (68) with the definition of $\mathcal{S}_\varpi$, we conclude that $\theta \cdot \mathbf{x}_\varpi + (1 - \theta) \cdot \mathbf{y}_\varpi \in \mathcal{S}_\varpi$ for all $\mathbf{x}_\varpi \in \mathcal{S}_\varpi$, $\mathbf{y}_\varpi \in \mathcal{S}_\varpi$, i.e., $\mathcal{S}_\varpi$ is convex.

To prove that $\mathcal{S}_\varpi$ is closed, we show that no point of $\overline{\mathcal{S}}_\varpi$ is a contact point of $\mathcal{S}_\varpi$. Indeed, the complement of $\mathcal{S}_\varpi$ in $\mathbb{R}^n$ is given by

$$\overline{\mathcal{S}}_\varpi = \big\{ \mathbf{y} \in \mathbb{R}^n : \textstyle\sum_{i \in (\mathcal{I}_n \setminus \mathcal{I}_n^\infty)} (\mathbf{Fy})_i^2 > 0 \big\}. \tag{69}$$

Let us consider some $\mathbf{y} \in \overline{\mathcal{S}}_\varpi$. Taking into account the definitions of $\mathcal{S}_\varpi$ and $\overline{\mathcal{S}}_\varpi$ and the fact that $\mathbf{F}$ is unitary, we have that, for any $\mathbf{x} \in \mathcal{S}_\varpi$,

$$\begin{aligned}
\|\mathbf{x} - \mathbf{y}\|_2^2 &= \|\mathbf{F}(\mathbf{x} - \mathbf{y})\|_2^2 = \|\mathbf{Fx} - \mathbf{Fy}\|_2^2 \\
&= \textstyle\sum_{i \in \mathcal{I}_n^\infty} ((\mathbf{Fx})_i - (\mathbf{Fy})_i)^2 + \sum_{i \in (\mathcal{I}_n \setminus \mathcal{I}_n^\infty)} (\mathbf{Fy})_i^2 \\
&\geqslant \textstyle\sum_{i \in (\mathcal{I}_n \setminus \mathcal{I}_n^\infty)} (\mathbf{Fy})_i^2 > 0.
\end{aligned} \tag{70}$$

Thus, for every $\mathbf{y} \in \overline{\mathcal{S}}_\varpi$, there exist no $\mathbf{x} \in \mathcal{S}_\varpi$ such that $\|\mathbf{x} - \mathbf{y}\|_2 < \epsilon$ with $\epsilon = \sqrt{\sum_{i \in (\mathcal{I}_n \setminus \mathcal{I}_n^\infty)} (\mathbf{Fy})_i^2}$, which means that none of the elements of $\overline{\mathcal{S}}_\varpi$ are contact points of $\mathcal{S}_\varpi$. Therefore, $\mathcal{S}_\varpi$ coincides with the set of its contact points, i.e., it is closed.

It follows from the definition of $\mathcal{S}_\varpi$ that $\mathbf{y} = (\mathbf{x} + \epsilon \cdot \mathbf{e}^{(1)}) \notin \mathcal{S}_\varpi$ for all $\mathbf{x} \in \mathcal{S}_\varpi$ and $\epsilon > 0$, where $\mathbf{e}^{(1)}$ is the unit vector with all but its first components equal to zero. Indeed, $(\mathbf{Fe}^{(1)})_i = 1/\sqrt{n} \neq 0$ for all $i \in \mathcal{I}_n$. Moreover, in that case, $\|\mathbf{x} - \mathbf{y}\|_2 = \epsilon$. Thus, for all $\mathbf{x} \in \mathcal{S}_\varpi$ and $\epsilon > 0$, there exists a $\mathbf{y} \in \mathbb{R}^n$ such that $\|\mathbf{x} - \mathbf{y}\|_2 < \epsilon$ and $\mathbf{y} \notin \mathcal{S}_\varpi$, which means that $\mathcal{S}_{\geqslant |\mathbf{s}|}$ has no interior points.

A necessary and sufficient condition for a subset $\mathcal{S}$ of a linear space $\mathbb{R}^n$ to be a subspace is that $(\alpha \cdot \mathbf{x} + \beta \cdot \mathbf{y}) \in \mathcal{S}$ for all $\alpha \in \mathbb{R}$, $\beta \in \mathbb{R}$, $\mathbf{x} \in \mathcal{S}$, and $\mathbf{y} \in \mathcal{S}$ (see, e.g., [2, p. 121]). That this applies to $\mathcal{S}_\varpi$ follows from the proof of its convexity above if we replace $\theta$ and $(1 - \theta)$ by, respectively, $\alpha$ and $\beta$.                                        ∎

**Proposition E.3.** *The intersection of sets $\mathcal{S}_{\geqslant |\mathbf{s}|}^\circ$ and $\mathcal{S}_\varpi$ is nonempty, i.e., there exists an $\mathbf{x} \in \mathcal{S}_{\geqslant |\mathbf{s}|}^\circ \cap \mathcal{S}_\varpi$. Consequently, $\mathcal{S}_{\geqslant |\mathbf{s}|} \cap \mathcal{S}_\varpi$ is also nonempty.*

*Proof.* Let us consider $\mathbf{x} = \lambda \cdot \mathbf{1}$, where $\mathbf{1}$ is an element of $\mathbb{R}^n$ with all its components equal to 1, and $\lambda > \max[\mathbf{s}]$. It follows directly from the definition of $\mathcal{S}_{\geqslant |\mathbf{s}|}$ that $\mathbf{x} \in \mathcal{S}_{\geqslant |\mathbf{s}|}^\circ \subset \mathcal{S}_{\geqslant |\mathbf{s}|}$. It also follows from the definitions of $\mathcal{S}_\varpi$ and unitary discrete Fourier transform that $(\mathbf{Fx})_i = \delta_{i,0} \cdot \sqrt{n} \cdot \lambda$, i.e., $\mathbf{x} \in \mathcal{S}_\varpi$. Therefore, $\mathbf{x} \in (\mathcal{S}_{\geqslant |\mathbf{s}|}^\circ \cap \mathcal{S}_\varpi) \subset (\mathcal{S}_{\geqslant |\mathbf{s}|} \cap \mathcal{S}_\varpi)$.                                        ∎

*Remark.* It is a direct consequence of *Propositions D.6*, *E.1*, and *E.2* that $\mathcal{S}_{\geqslant |\mathbf{s}|} \cap \mathcal{S}_\varpi$ is also closed and convex.

*Metric projections onto $\mathcal{S}_{\geqslant |\mathbf{s}|}$ and $\mathcal{S}_\varpi$*

The metric projections of any point in $\mathbb{R}^n$ onto its convex subsets relevant to us, i.e., $\mathcal{S}_{\geqslant |\mathbf{s}|}$ and $\mathcal{S}_\varpi$, are achieved by the following operators.

**Proposition E.4.** *The metric projection operator* $\mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}}$ *is defined by the elementwise maximum of the target signal* $\mathbf{s}$ *and the input argument* $\mathbf{z}$:

$$\mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}}[\mathbf{z}] = |\mathbf{s}| + (\mathbf{z} - |\mathbf{s}|) \circ \theta(\mathbf{z} - |\mathbf{s}|). \tag{71}$$

*Proof.* Note that

$$\left(\mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}}[\mathbf{z}]\right)_i = \begin{cases} z_i, & \text{if } z_i \geqslant |s_i| \\ |s_i|, & \text{if } z_i < |s_i| \end{cases} \qquad \forall i \in \mathcal{I}_n. \tag{72}$$

Hence, a necessary and sufficient condition for transforming any $\mathbf{z} \in \mathbb{R}^n$ to $\mathbf{x} \in \mathcal{S}_{\geqslant|\mathbf{s}|}$ is to increase every component $z_i$ of $\mathbf{z}$ that does not satisfy $z_i \geqslant |s_i|$ by at least $|s_i| - z_i$, independent of values of the remaining components. Now, we have from the definition of the Euclidean norm that

$$\|\mathbf{z} - \mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n (z_i - x_i)^2} \qquad \forall \mathbf{z} \in \mathbb{R}^n, \forall \mathbf{x} \in \mathcal{S}_{\geqslant|\mathbf{s}|}. \tag{73}$$

Thus, for every $\mathbf{z} \in \mathbb{R}^n$, $\|\mathbf{z}-\mathbf{x}\|_2$ is minimized by an $\mathbf{x} \in \mathcal{S}_{\geqslant|\mathbf{s}|}$ that is obtained by incrementing all components of $\mathbf{z}$ that satisfy $z_i < |s_i|$ by no more than necessary, i.e., by $|s_i| - z_i$, and leaving the remaining components intact. However, this is precisely how the operator $\mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}}$ is defined via (72). Therefore, by using the *Definition D.13* of the projection, we conclude that $\mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}}$ projects $\mathbf{z} \in \mathbb{R}^n$ onto $\mathcal{S}_{\geqslant|\mathbf{s}|}$. ∎

**Proposition E.5.** *The metric projection operator* $\mathbf{P}_{\mathcal{S}_{\varpi}}$ *is defined by a rectangular low-pass-filter transformation*

$$\mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{z}] = (\mathbf{F}^{-1}\,\mathbf{W}_{\varpi}\,\mathbf{F})\,\mathbf{z} = \sum_{i \in \mathcal{I}_n^{\varpi}} \langle \mathbf{f}^{(i)}, \mathbf{z} \rangle \cdot \mathbf{f}^{(i)}. \tag{74}$$

*Here,* $\mathbf{f}^{(i)}$ *is the $i$-th column of the* DFT *matrix* $\mathbf{F}$. $\mathbf{W}_{\varpi}$ *is a diagonal matrix such that*

$$(W_{\varpi})_{ii} = \begin{cases} 1, & \text{if } i \in \mathcal{I}_n^{\varpi} \\ 0, & \text{otherwise} \end{cases}. \tag{75}$$

*Proof.* Let us consider an element of the set $\mathcal{S}_{\varpi}$ expressed by $\mathbf{x} = \sum_{i \in \mathcal{I}_n^{\varpi}} a_i \cdot \mathbf{f}^{(i)}$. It follows from the definition of the Euclidean norm (see Section II in the main text) that

$$\begin{aligned}
\|\mathbf{x} - \mathbf{z}\|_2^2 &= \left\langle \left(\mathbf{z} - \sum_{i \in \mathcal{I}_n^{\varpi}} a_i \cdot \mathbf{f}^{(i)}\right), \left(\mathbf{z} - \sum_{i \in \mathcal{I}_n^{\varpi}} a_i \cdot \mathbf{f}^{(i)}\right) \right\rangle \\
&= \langle \mathbf{z}, \mathbf{z} \rangle - 2 \cdot \sum_{i \in \mathcal{I}_n^{\varpi}} a_i \cdot \langle \mathbf{f}^{(i)}, \mathbf{z} \rangle + \sum_{i \in \mathcal{I}_n^{\varpi}} a_i^2 \\
&= \langle \mathbf{z}, \mathbf{z} \rangle - \sum_{i \in \mathcal{I}_n^{\varpi}} \langle \mathbf{f}^{(i)}, \mathbf{z} \rangle^2 + \sum_{i \in \mathcal{I}_n^{\varpi}} \left(a_i - \langle \mathbf{f}^{(i)}, \mathbf{z} \rangle\right)^2.
\end{aligned} \tag{76}$$

When writing the second equality above, we used the fact that $\{\mathbf{f}^{(1)}, \mathbf{f}^{(2)}, \ldots, \mathbf{f}^{(n)}\}$ are orthonormal. It follows from the last equality of (76) that $\|\mathbf{x} - \mathbf{z}\|_2$, as a function of $a_i$, is minimized by $a_i = \langle \mathbf{f}^{(i)}, \mathbf{z} \rangle$ for every $i \in \mathcal{I}_n^{\varpi}$. Thus, by *Definition D.13*, $\sum_{i \in \mathcal{I}_n^{\varpi}} \langle \mathbf{f}^{(i)}, \mathbf{z} \rangle \cdot \mathbf{f}^{(i)} = \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{z}]$ is the projection of $\mathbf{z}$ onto $\mathcal{S}_{\varpi}$. ∎

*Remark.* $\mathbf{P}_{\mathcal{S}_{\varpi}}$ is a linear operator, whereas $\mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}}$ is not.

## F. CONVERGENCE PROOFS

Here, we provide proofs of the propositions concerning the convergence of the AP algorithms that are formulated in Section III of the main text. The proofs are adapted for finite-dimensional Euclidean spaces and exploit the particular structure of the modulator constraint sets.[15] For the sake of convenience, we repeat the original assertions as well.

### *AP-B algorithm*

**Proposition III.1.** *A sequence* $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ *formed by the AP-B algorithm for* $\epsilon_{tol} = 0$ *and* $N_{iter} \to +\infty$ *converges to some* $\mathbf{m}^{\dagger} \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$. *The convergence is geometric and monotonic, i.e., there exist* $\gamma > 0$ *and* $0 < r < 1$ *such that* $\|\mathbf{m}^{(i)} - \mathbf{m}^{\dagger}\|_2 \leqslant \gamma \cdot r^i$ *and* $\|\mathbf{m}^{(i+1)} - \mathbf{m}^{\dagger}\|_2 \leqslant \|\mathbf{m}^{(i)} - \mathbf{m}^{\dagger}\|_2$ *for* $i \geqslant 0$.

*Proof.* If the sequence $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ terminates with some $\mathbf{m}^{(N)}$, i.e., $\mathbf{m}^{(N+j)} = \mathbf{m}^{(N)}$ for every $j > 0$, it follows from the formulation of the AP-B algorithm that $\mathbf{m}^{(N)} = \mathbf{a}^{(N)}$, i.e., $\mathbf{m}^{(N)} \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$.[16] Thus, $\mathbf{m}^{(N)} = \mathbf{m}^{\dagger}$, which means that the solution is achieved in a finite number of iterations. We next consider the case when the sequence $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ is infinite. The rest of the proof is divided into three parts for clarity.

Convergence. The outline of the convergence proof is as follows. We first demonstrate that the distance between any $\mathbf{x} \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$ and $\mathbf{m}^{(i)}$ or $\mathbf{a}^{(i)}$ decreases with every iteration. Using this, we then show that the sequence $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ is bounded, and thus, by the Bolzano-Weierstrass theorem and closedness of $\mathcal{S}_{\geqslant|\mathbf{s}|}$ and $\mathcal{S}_{\varpi}$, has a subsequence that converges to some $\mathbf{m}^{\dagger} \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$. Referring to the first result again (that the distance between any $\mathbf{x} \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$ and $\mathbf{m}^{(i)}$ decreases with every iteration), we finally deduce that the original sequence $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ converges to the same $\mathbf{m}^{\dagger} \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$ as any of its infinite subsequences.

$\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$ is nonempty by *Proposition E.3*. Let us consider some $\mathbf{x} \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$ together with $\mathbf{m}^{(i)}$ and $\mathbf{a}^{(i)}$ taken from the sequences $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ and $\mathbf{a}^{(0)}, \mathbf{a}^{(1)}, \ldots, \mathbf{a}^{(i)}, \ldots$ for some $i \geqslant 0$. Then, we have

$$\|\mathbf{x} - \mathbf{m}^{(i)}\|_2^2 = \|\mathbf{x} - \mathbf{a}^{(i+1)} + \mathbf{a}^{(i+1)} - \mathbf{m}^{(i)}\|_2^2$$
$$= \underbrace{\|\mathbf{x} - \mathbf{a}^{(i+1)}\|_2^2}_{\geqslant 0} + \underbrace{\|\mathbf{a}^{(i+1)} - \mathbf{m}^{(i)}\|_2^2}_{\geqslant 0} + 2 \cdot \underbrace{\left\langle \mathbf{m}^{(i)} - \mathbf{a}^{(i+1)}, \mathbf{a}^{(i+1)} - \mathbf{x} \right\rangle}_{\geqslant 0}. \tag{77}$$

The nonnegativity of the last term in the second line of (77) follows from the containing-half-space inequality (65) and $\mathbf{a}^{(i+1)} = \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{m}^{(i)}]$. Therefore, (77) implies that

$$\|\mathbf{x} - \mathbf{m}^{(i)}\|_2^2 \geqslant \|\mathbf{x} - \mathbf{a}^{(i+1)}\|_2^2 + \|\mathbf{a}^{(i+1)} - \mathbf{m}^{(i)}\|_2^2 \qquad \forall i \geqslant 0 \tag{78}$$

---

[15]For the foundations of AP algorithms in a more general context of arbitrary closed convex subsets of Hilbert spaces, we refer an interested reader to the seminal works by Bregman [4], Gurin et al. [5], and Boyle & Dykstra [6].

[16]We remind the reader that $\mathbf{a}^{(i)} = \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{m}^{(i-1)}]$ for any $i > 0$ in the case of the AP-B algorithm.

and

$$\|\mathbf{x} - \mathbf{m}^{(i)}\|_2 \geqslant \|\mathbf{x} - \mathbf{a}^{(i+1)}\|_2 \qquad \forall i \geqslant 0. \tag{79}$$

Replacing $\mathbf{m}^{(i)}$ by $\mathbf{a}^{(i+1)}$ and $\mathbf{a}^{(i+1)}$ by $\mathbf{m}^{(i+1)}$ in (77), and using the same argumentation as above, including $\mathbf{m}^{(i+1)} = \mathbf{P}_{\mathcal{S}_{\geqslant |\mathbf{S}|}}[\mathbf{a}^{(i+1)}]$, we deduce that

$$\|\mathbf{x} - \mathbf{a}^{(i+1)}\|_2^2 \geqslant \|\mathbf{x} - \mathbf{m}^{(i+1)}\|_2^2 + \|\mathbf{m}^{(i+1)} - \mathbf{a}^{(i+1)}\|_2^2 \qquad \forall i \geqslant -1 \tag{80}$$

and

$$\|\mathbf{x} - \mathbf{a}^{(i+1)}\|_2 \geqslant \|\mathbf{x} - \mathbf{m}^{(i+1)}\|_2 \qquad \forall i \geqslant -1. \tag{81}$$

The validity of (80) and (81) for not only $i \geqslant 0$ but also $i = -1$ follows from the particular initial conditions of the AP-B algorithm. Combining (79) and (81) yields

$$\|\mathbf{x} - \mathbf{m}^{(0)}\|_2 \geqslant \|\mathbf{x} - \mathbf{m}^{(1)}\|_2 \geqslant \ldots \geqslant \|\mathbf{x} - \mathbf{m}^{(i)}\|_2 \geqslant \ldots \tag{82}$$

and, equivalently,

$$\|\mathbf{x} - \mathbf{m}^{(0)}\|_2^2 \geqslant \|\mathbf{x} - \mathbf{m}^{(1)}\|_2^2 \geqslant \ldots \geqslant \|\mathbf{x} - \mathbf{m}^{(i)}\|_2^2 \geqslant \ldots. \tag{83}$$

(83) states that the sequence $\|\mathbf{x} - \mathbf{m}^{(0)}\|_2^2, \|\mathbf{x} - \mathbf{m}^{(1)}\|_2^2, \ldots \|\mathbf{x} - \mathbf{m}^{(i)}\|_2^2, \ldots$ is monotonically decreasing. This sequence is bounded from below by 0 because of the nonnegativity of the norm, and therefore, it converges to its greatest lower bound $L \geqslant 0$ by the monotone convergence theorem (see *Proposition D.10*). Thus, for every $\epsilon > 0$, there exists an $N(\epsilon)$ such that $0 \leqslant \|\mathbf{x} - \mathbf{m}^{(i)}\|_2^2 - L \leqslant \epsilon$ whenever $i > N(\epsilon)$. It follows then from (79) and (81) that $L \leqslant \|\mathbf{x} - \mathbf{a}^{(i+1)}\|_2^2 \leqslant \|\mathbf{x} - \mathbf{m}^{(i)}\|_2^2 \leqslant L + \epsilon$, so that $0 \leqslant \|\mathbf{x} - \mathbf{m}^{(i)}\|_2^2 - \|\mathbf{x} - \mathbf{a}^{(i+1)}\|_2^2 < \epsilon$, and because of (78), also $0 \leqslant \|\mathbf{a}^{(i+1)} - \mathbf{m}^{(i)}\|_2 < \sqrt{\epsilon}$ whenever $i > N(\epsilon)$, i.e., the sequence $\|\mathbf{a}^{(1)} - \mathbf{m}^{(0)}\|_2, \|\mathbf{a}^{(2)} - \mathbf{m}^{(1)}\|_2, \ldots, \|\mathbf{a}^{(i)} - \mathbf{m}^{(i-1)}\|_2, \ldots$ converges to 0. If the sequence converges, then any of its infinite subsequences (see *Definition D.11*) converges as well, because the removal of elements from the sequence does not change the validity of the convergence condition:

$$\forall \epsilon > 0 \ \ \exists N'(\epsilon) : \ i > N'(\epsilon) \implies \|\mathbf{a}^{(k_i+1)} - \mathbf{m}^{(k_i)}\|_2 < \epsilon, \tag{84}$$

where $k_i > k_j$ for $i > j \geqslant 0$ and $k_i \geqslant i$ for $i \geqslant 0$.

Next, we have from the triangle inequality (63) that

$$\|\mathbf{m}^{(i)} - \mathbf{m}^{(j)}\|_2 \leqslant \|\mathbf{x} - \mathbf{m}^{(i)}\|_2 + \|\mathbf{x} - \mathbf{m}^{(j)}\|_2 \qquad \forall i, j \geqslant 0. \tag{85}$$

Moreover, according to (82), $\|\mathbf{x} - \mathbf{m}^{(i)}\|_2 \leqslant \|\mathbf{x} - \mathbf{m}^{(0)}\|_2$ for $i \geqslant 0$. Thus, for all $i, j \geqslant 0$,

$$\|\mathbf{m}^{(i)} - \mathbf{m}^{(j)}\|_2 \leqslant 2 \cdot \|\mathbf{x} - \mathbf{m}^{(0)}\|_2, \tag{86}$$

i.e., the sequence $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ is bounded (see *Definition D.7*). Consequently, according to the Bolzano-Weierstrass theorem (see *Proposition D.12*), this sequence has a subsequence $\mathbf{m}^{(k_0)}, \mathbf{m}^{(k_1)}, \ldots, \mathbf{m}^{(k_i)}, \ldots$ that converges to some $\mathbf{m}^\dagger \in \mathbb{R}^n$:

$$\forall \epsilon > 0 \ \ \exists N''(\epsilon): \ i > N''(\epsilon) \implies \|\mathbf{m}^\dagger - \mathbf{m}^{(k_i)}\|_2 < \epsilon. \tag{87}$$

We show now that $\mathbf{m}^\dagger \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi$. By construction, $\mathbf{m}^{(i)} \in \mathcal{S}_{\geqslant|\mathbf{s}|}$ for every $i \geqslant 0$. According to (87), there exists an $\mathbf{m}^{(i)}$ for any $\epsilon > 0$ such that $\|\mathbf{m}^{(i)} - \mathbf{m}^\dagger\|_2 < \epsilon$. Hence, $\mathbf{m}^\dagger$ is a contact point of $\mathcal{S}_{\geqslant|\mathbf{s}|}$ (see *Definition D.4*). Moreover, because the latter set is closed (see *Definition D.5* and *Proposition E.1*), $\mathbf{m}^\dagger \in \mathcal{S}_{\geqslant|\mathbf{s}|}$. Next, by exploiting the triangle inequality (63), we can write

$$\|\mathbf{m}^\dagger - \mathbf{a}^{(k_i+1)}\|_2 \leqslant \|\mathbf{a}^{(k_i+1)} - \mathbf{m}^{(k_i)}\|_2 + \|\mathbf{m}^\dagger - \mathbf{m}^{(k_i)}\|_2 \qquad \forall i \geqslant 0. \tag{88}$$

Combining (88) with (84) and (87) and introducing $N'''(\epsilon) = \max[N'(\epsilon/2), N''(\epsilon/2)]$, we get

$$\forall \epsilon > 0 \ \ \exists N'''(\epsilon): \ i > N'''(\epsilon) \implies \|\mathbf{m}^\dagger - \mathbf{a}^{(k_i+1)}\|_2 < \epsilon, \tag{89}$$

i.e., the subsequence $\mathbf{a}^{(k_0+1)}, \mathbf{a}^{(k_1+1)}, \ldots, \mathbf{a}^{(k_i+1)}, \ldots$ converges to $\mathbf{m}^\dagger$. The set $\mathcal{S}_\varpi$ is closed (see *Proposition E.2*), and $\mathbf{a}^{(i)} \in \mathcal{S}_\varpi$ for every $i \geqslant 0$ by construction. Therefore, using the same argumentation as for the subsequence $\mathbf{m}^{(k_0)}, \mathbf{m}^{(k_1)}, \ldots, \mathbf{m}^{(k_i)}, \ldots$, we conclude that $\mathbf{m}^\dagger \in \mathcal{S}_\varpi$.

Finally, because $\mathbf{m}^\dagger \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi$, (82) gives

$$\|\mathbf{m}^\dagger - \mathbf{m}^{(0)}\|_2 \geqslant \|\mathbf{m}^\dagger - \mathbf{m}^{(1)}\|_2 \geqslant \ldots \geqslant \|\mathbf{m}^\dagger - \mathbf{m}^{(i)}\|_2 \geqslant \ldots. \tag{90}$$

In the light of (90), the statement of (87) generalizes to

$$\forall \epsilon > 0 \ \ \exists N''''(\epsilon): \ i > N''''(\epsilon) \implies \|\mathbf{m}^\dagger - \mathbf{m}^{(i)}\|_2 < \epsilon, \tag{91}$$

where $N''''(\epsilon) = k_{N''(\epsilon)+1}$. Thus, the sequence $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ converges to $\mathbf{m}^\dagger \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi$.

Monotonicity. The monotonicity of the convergence of the sequence $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ to $\mathbf{m}^\dagger$ is declared by (90).

Rate. The key point in establishing the geometric convergence of the AP-B algorithm is the fact that the intersection of $\mathcal{S}_\varpi$ and the interior of $\mathcal{S}_{\geqslant|\mathbf{s}|}$ is nonempty. Using this fact, we first show that the distances between $\mathbf{m}^{(i)}$ and $\mathcal{S}_{\geqslant|\mathbf{s}|}^\circ \cap \mathcal{S}_\varpi$ or $\mathbf{a}^{(i+1)}$ and $\mathcal{S}_{\geqslant|\mathbf{s}|}^\circ \cap \mathcal{S}_\varpi$ can be bounded by, respectively, $\|\mathbf{m}^{(i)} - \mathbf{a}^{(i+1)}\|_2$ or $\|\mathbf{m}^{(i+1)} - \mathbf{a}^{(i+1)}\|_2$ multiplied by a universal factor that is greater than one and independent of the iteration number $i$. In the next step, we exploit properties of metric projections to obtain two additional inequalities, which, in combination with the first result, allow deriving a decreasing geometric sequence that bounds $\|\mathbf{m}^{(0)} - \mathbf{m}^\dagger\|_2, \|\mathbf{m}^{(1)} - \mathbf{m}^\dagger\|_2, \ldots, \|\mathbf{m}^{(i)} - \mathbf{m}^\dagger\|_2, \ldots$ from above.

To start, let us consider some $\mathbf{x} \in \mathcal{S}_{\geq|\mathbf{s}|}^{\circ} \cap \mathcal{S}_{\varpi}$. Such an element exists according to *Proposition E.3*. Also, there exists an $\epsilon > 0$ such that $\mathbf{y} \in \mathcal{S}_{\geq|\mathbf{s}|}$ if $\|\mathbf{x} - \mathbf{y}\|_2 < \epsilon$ because $\mathbf{x}$ belongs to the interior of $\mathcal{S}_{\geq|\mathbf{s}|}$ (see *Definition D.2*). If so, then it is also possible to choose a positive $\beta < \epsilon$ such that $\mathbf{y} \in \mathcal{S}_{\geq|\mathbf{s}|}$ if $\|\mathbf{x} - \mathbf{y}\|_2 \leq \beta$. We now introduce

$$\mathbf{z}^{(i)} = \frac{\alpha_i}{\alpha_i + \beta} \cdot \mathbf{x} + \frac{\beta}{\alpha_i + \beta} \cdot \mathbf{a}^{(i)}, \quad i \geq 0, \tag{92}$$

where $\alpha_i = \|\mathbf{a}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geq|\mathbf{s}|}}[\mathbf{a}^{(i)}]\|_2 = \|\mathbf{a}^{(i)} - \mathbf{m}^{(i)}\|_2$. Note that $\alpha_i/(\alpha_i+\beta) \in (0,1)$, and $\beta/(\alpha_i+\beta) = 1-\alpha_i/(\alpha_i+\beta)$. Moreover, $\mathbf{x} \in \mathcal{S}_{\varpi}$, and $\mathbf{a}^{(i)} \in \mathcal{S}_{\varpi}$ by construction. Therefore, by the *Definition D.1* of a convex set and the fact that $\mathcal{S}_{\varpi}$ is convex (see *Proposition E.2*), we have $\mathbf{z}^{(i)} \in \mathcal{S}_{\varpi}$. On the other hand, (92) can be rewritten as

$$\mathbf{z}^{(i)} = \frac{\alpha_i}{\alpha_i + \beta} \cdot \underbrace{\left(\mathbf{x} + \frac{\beta}{\alpha_i} \cdot (\mathbf{a}^{(i)} - \mathbf{m}^{(i)})\right)}_{\mathbf{y}'} + \frac{\beta}{\alpha_i + \beta} \cdot \mathbf{m}^{(i)}. \tag{93}$$

In the above expression, $\|\mathbf{x} - \mathbf{y}'\|_2 = \beta$. Therefore, $\mathbf{y}' \in \mathcal{S}_{\geq|\mathbf{s}|}$ by the definition of $\beta$. Moreover, $\mathbf{m}^{(i)} \in \mathcal{S}_{\geq|\mathbf{s}|}$ by construction, which implies $\mathbf{z}^{(i)} \in \mathcal{S}_{\geq|\mathbf{s}|}$ because $\mathcal{S}_{\geq|\mathbf{s}|}$ is convex (see *Proposition E.1*). Hence, altogether, we conclude that $\mathbf{z}^{(i)} \in \mathcal{S}_{\geq|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$ for $i \geq 0$.

Based on the above consideration, we show now that

$$\|\mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geq|\mathbf{s}|} \cap \mathcal{S}_{\varpi}}[\mathbf{m}^{(i)}]\|_2 \leq \|\mathbf{m}^{(i)} - \mathbf{a}^{(i+1)}\|_2 \cdot (1 + \|\mathbf{x} - \mathbf{m}^{(0)}\|_2/\beta) \tag{94}$$

and

$$\|\mathbf{a}^{(i+1)} - \mathbf{P}_{\mathcal{S}_{\geq|\mathbf{s}|} \cap \mathcal{S}_{\varpi}}[\mathbf{a}^{(i+1)}]\|_2 \leq \|\mathbf{m}^{(i+1)} - \mathbf{a}^{(i+1)}\|_2 \cdot (1 + \|\mathbf{x} - \mathbf{m}^{(0)}\|_2/\beta) \tag{95}$$

for $i \geq 0$. To demonstrate (94), note that

$$\begin{aligned}
\|\mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geq|\mathbf{s}|} \cap \mathcal{S}_{\varpi}}[\mathbf{m}^{(i)}]\|_2 &\leq \|\mathbf{m}^{(i)} - \mathbf{z}^{(i)}\|_2 \\
&\leq \|\mathbf{m}^{(i)} - \mathbf{a}^{(i+1)}\|_2 + \|\mathbf{a}^{(i+1)} - \mathbf{z}^{(i+1)}\|_2 \\
&= \|\mathbf{m}^{(i)} - \mathbf{a}^{(i+1)}\|_2 + \frac{\alpha_{i+1}}{\alpha_{i+1} + \beta} \cdot \|\mathbf{x} - \mathbf{a}^{(i+1)}\|_2 \\
&\leq \|\mathbf{m}^{(i)} - \mathbf{a}^{(i+1)}\|_2 + \frac{\alpha_{i+1}}{\beta} \cdot \|\mathbf{x} - \mathbf{m}^{(0)}\|_2 \\
&= \|\mathbf{m}^{(i)} - \mathbf{a}^{(i+1)}\|_2 + \frac{\|\mathbf{m}^{(i+1)} - \mathbf{a}^{(i+1)}\|_2}{\beta} \cdot \|\mathbf{x} - \mathbf{m}^{(0)}\|_2 \\
&\leq \|\mathbf{m}^{(i)} - \mathbf{a}^{(i+1)}\|_2 + \frac{\|\mathbf{m}^{(i)} - \mathbf{a}^{(i+1)}\|_2}{\beta} \cdot \|\mathbf{x} - \mathbf{m}^{(0)}\|_2.
\end{aligned} \tag{96}$$

In (96), we used the fact that $\mathbf{z} \in \mathcal{S}_{\geq|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$ and the *Definition D.13* of the projection operator (the first line), the triangle inequality (63) (the second line), (92) (the third line), combined inequalities (79) and (81) (the fourth

line), and the *Definition D.13* of the projection operator again (the last line). Similarly to (96), we can write

$$
\begin{aligned}
\|\mathbf{a}^{(i+1)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{a}^{(i+1)}]\|_2 &\leqslant \|\mathbf{a}^{(i+1)} - \mathbf{z}^{(i+1)}\|_2 \\
&= \frac{\alpha_{i+1}}{\alpha_{i+1}+\beta}\cdot\|\mathbf{x}-\mathbf{a}^{(i+1)}\|_2 \\
&\leqslant \frac{\alpha_{i+1}}{\beta}\cdot\|\mathbf{x}-\mathbf{m}^{(0)}\|_2 \\
&\leqslant \alpha_{i+1} + \frac{\alpha_{i+1}}{\beta}\cdot\|\mathbf{x}-\mathbf{m}^{(0)}\|_2 \\
&= \|\mathbf{m}^{(i+1)}-\mathbf{a}^{(i+1)}\|_2 + \frac{\|\mathbf{m}^{(i+1)}-\mathbf{a}^{(i+1)}\|_2}{\beta}\cdot\|\mathbf{x}-\mathbf{m}^{(0)}\|_2,
\end{aligned}
\tag{97}
$$

which proves (95).

Next, we derive two additional inequalities. In particular, we have

$$
\|\mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2^2 - \|\mathbf{a}^{(i+1)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{a}^{(i+1)}]\|_2^2
$$

$$
\begin{aligned}
&\geqslant \|\mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2^2 - \|\mathbf{a}^{(i+1)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2^2 \\
&= \|\mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2^2 - \|(\mathbf{a}^{(i+1)}-\mathbf{m}^{(i)}) + (\mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{m}^{(i)}])\|_2^2 \\
&= -\|\mathbf{a}^{(i+1)}-\mathbf{m}^{(i)}\|_2^2 + 2\cdot\langle\mathbf{a}^{(i+1)}-\mathbf{m}^{(i)}, \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{m}^{(i)}]-\mathbf{m}^{(i)}\rangle \\
&= -\|\mathbf{a}^{(i+1)}-\mathbf{m}^{(i)}\|_2^2 + 2\cdot\langle\mathbf{a}^{(i+1)}-\mathbf{m}^{(i)}, (\mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{m}^{(i)}]-\mathbf{a}^{(i+1)}) + (\mathbf{a}^{(i+1)}-\mathbf{m}^{(i)})\rangle \\
&= \|\mathbf{a}^{(i+1)}-\mathbf{m}^{(i)}\|_2^2 + 2\cdot\underbrace{\langle\mathbf{m}^{(i)}-\mathbf{a}^{(i+1)}, \mathbf{a}^{(i+1)}-\mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\rangle}_{\geqslant 0} \\
&\geqslant \|\mathbf{a}^{(i+1)}-\mathbf{m}^{(i)}\|_2^2
\end{aligned}
\tag{98}
$$

for $i \geqslant 0$. Thus,

$$
\|\mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2^2 - \|\mathbf{a}^{(i+1)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{a}^{(i+1)}]\|_2^2 \geqslant \|\mathbf{m}^{(i)}-\mathbf{a}^{(i+1)}\|_2^2, \quad i \geqslant 0.
\tag{99}
$$

In (98), we used the *Definition D.13* of the projection operator (the second line) and the containing-half-space inequality (65) along with $\mathbf{m}^{(i)} = \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}}[\mathbf{a}^{(i)}]$ (the sixth line). Replacing $\mathbf{a}^{(i+1)}$ by $\mathbf{m}^{(i+1)}$ and $\mathbf{m}^{(i)}$ by $\mathbf{a}^{(i+1)}$ and repeating the same steps as in (98), we obtain

$$
\|\mathbf{a}^{(i+1)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{a}^{(i+1)}]\|_2^2 - \|\mathbf{m}^{(i+1)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{m}^{(i+1)}]\|_2^2 \geqslant \|\mathbf{m}^{(i+1)}-\mathbf{a}^{(i+1)}\|_2^2, \quad i \geqslant 0.
\tag{100}
$$

Combining (94) with (99) and (95) with (100), we obtain, respectively,

$$
\left(1 - \frac{1}{(1+\|\mathbf{x}-\mathbf{m}^{(0)}\|_2/\beta)^2}\right)\cdot\|\mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2^2 \geqslant \|\mathbf{a}^{(i+1)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{a}^{(i+1)}]\|_2^2
\tag{101}
$$

and

$$
\left(1 - \frac{1}{(1+\|\mathbf{x}-\mathbf{m}^{(0)}\|_2/\beta)^2}\right)\cdot\|\mathbf{a}^{(i+1)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{a}^{(i+1)}]\|_2^2 \geqslant \|\mathbf{m}^{(i+1)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}\cap\mathcal{S}_\varpi}[\mathbf{m}^{(i+1)}]\|_2^2,
\tag{102}
$$

which then lead to

$$\underbrace{\left(1 - \frac{1}{(1 + \|\mathbf{x} - \mathbf{m}^{(0)}\|_2/\beta)^2}\right)}_{r<1} \cdot \|\mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2 \geqslant \|\mathbf{m}^{(i+1)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i+1)}]\|_2 \tag{103}$$

for $i \geqslant 0$. Starting with $i = 0$ and applying (103) iteratively, we get

$$r^i \cdot \|\mathbf{m}^{(0)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(0)}]\|_2 \geqslant \|\mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2, \quad i \geqslant 0. \tag{104}$$

According to the triangle inequality (63),

$$\|\mathbf{m}^{(i)} - \mathbf{m}^\dagger\|_2 \leqslant \|\mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2 + \|\mathbf{m}^\dagger - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2 \tag{105}$$

and [see (64)]

$$\|\mathbf{m}^{(j)} - \mathbf{m}^\dagger\|_2 \geqslant \left|\|\mathbf{m}^{(j)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2 - \|\mathbf{m}^\dagger - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2\right|. \tag{106}$$

(106) and (91) together imply that a sequence

$$\|\mathbf{m}^{(0)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2, \|\mathbf{m}^{(1)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2, \ldots, \|\mathbf{m}^{(j)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2, \ldots \tag{107}$$

converges to $\|\mathbf{m}^\dagger - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2$ for every $i \geqslant 0$. On the other hand, this sequence is monotonically decreasing [see (82)], and thus, by the monotone convergence theorem (*Proposition D.10*), it converges to its greatest lower bound, i.e., $\|\mathbf{m}^{(j)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2 \geqslant \|\mathbf{m}^\dagger - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2$ for all $i, j \geqslant 0$. Consequently, (105) reduces to

$$\|\mathbf{m}^{(i)} - \mathbf{m}^\dagger\|_2 \leqslant 2 \cdot \|\mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2. \tag{108}$$

Applying (108) to (104), we finally obtain

$$\gamma \cdot r^i \geqslant \|\mathbf{m}^{(i)} - \mathbf{m}^\dagger\|_2, \quad i \geqslant 0, \tag{109}$$

where $\gamma = 2 \cdot \|\mathbf{m}^{(0)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi}[\mathbf{m}^{(0)}]\|_2 > 0$.                                   ∎

*Remark.* 1) The convergence proof of the iterative scheme AP-B relies entirely on the convexity and closedness of the constraint sets. Therefore, this algorithm extends to more general sets than $\mathcal{S}_{\geqslant|\mathbf{s}|}$ and $\mathcal{S}_\varpi$. 2) The geometric nature of the convergence requires additionally that the interior of at least one of the constraint sets is nonempty and shares elements with the other set.

**Proposition III.2.** *Consider* $\mathbf{m} \in \mathcal{M}_\omega$ *and* $\mathbf{c} \in \mathcal{C}_d$ *with* $|c_j| = \sum_{k=1}^{n/\nu} (\tilde{c}_{\nu \cdot k} \cdot e^{\imath 2\pi \nu (k-1)(j-1)/n})$, *where* $\tilde{c}_{\nu \cdot k} \in \mathbb{C}$ *and* $n/\nu \in \mathbb{N}$. *If* $\varpi \geqslant \omega$ *and* $\nu \geqslant \varpi + \omega - 1$, *then a sequence* $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ *formed by the AP-B algorithm for* $\epsilon_{tol} = 0$ *and* $N_{iter} \to +\infty$ *converges to* $\mathbf{m}$.

*Proof.* The proof relies on two auxiliary results that apply to the $\mathbf{m}$ and $\mathbf{c}$ specified in the proposition:

- For every $q \in \mathbb{R}$,

$$\mathbf{z} = q + (|\mathbf{c}| - q) \circ \theta(|\mathbf{c}| - q) \implies z_j = \sum_{k=1}^{n/\nu} \left( \tilde{z}_{\nu \cdot k} \cdot e^{\imath 2\pi\nu(k-1)(j-1)/n} \right), \quad j \in \mathcal{I}_n. \tag{110}$$

- For $\mathbf{z}$ defined by (110),

$$\mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m} \circ \mathbf{z}] = \langle \mathbf{z} \rangle \cdot \mathbf{m}, \tag{111}$$

where, $\langle |\mathbf{z}| \rangle = \frac{1}{n} \sum_{i=1}^{n} |z_i|$.

To show (110), consider a $\mathbf{g} \in \mathbb{R}^n$ whose elements form a periodic sequence with the fundamental frequency $\nu \in \mathbb{N}$ such that $n/\nu \in \mathbb{N}$. Like any element of $\mathbb{R}^n$, $\mathbf{g}$ can be expressed through its DFT:

$$g_j = \frac{1}{\sqrt{n}} \sum_{k=1}^{n} \left( \tilde{g}_k \cdot e^{\imath 2\pi(k-1)(j-1)/n} \right), \quad j \in \mathcal{I}_n. \tag{112}$$

On the other hand, the periodicity of $g_1, g_2, \ldots, g_n$ implies that, for every $k \in \mathcal{I}_n$,

$$\begin{aligned}
\tilde{g}_k = (\mathbf{F}\mathbf{g})_k &= \frac{1}{\sqrt{n}} \sum_{j=1}^{n} \left( g_j \cdot e^{-\imath 2\pi(k-1)(j-1)/n} \right) \\
&= \frac{1}{\sqrt{n}} \sum_{j=1}^{n/\nu} g_j \sum_{l=1}^{\nu} e^{-\imath 2\pi((l-1)\cdot(n/\nu)+(j-1))(k-1)/n} \\
&= \frac{1}{\sqrt{n}} \sum_{j=1}^{n/\nu} \left( g_j \cdot e^{-\imath 2\pi(k-1)(j-1)/n} \sum_{l=1}^{\nu} \left( e^{-\imath 2\pi(k-1)/\nu} \right)^{l-1} \right) \\
&= \underbrace{\left( \frac{1 - e^{-\imath 2\pi(k-1)}}{1 - e^{-\imath 2\pi(k-1)/\nu}} \right)}_{=0, \text{ if } ((k-1)/\nu) \notin \mathbb{N}} \cdot \frac{1}{\sqrt{n}} \sum_{j=1}^{n/\nu} \left( g_j \cdot e^{-\imath 2\pi(k-1)(j-1)/n} \right).
\end{aligned} \tag{113}$$

Combining (112) and (113) gives

$$g_j = \frac{1}{\sqrt{n}} \sum_{k=1}^{n/\nu} \left( \tilde{g}_{\nu \cdot k} \cdot e^{\imath 2\pi\nu(k-1)(j-1)/n} \right), \quad j \in \mathcal{I}_n, \tag{114}$$

where $\tilde{g}_{\nu \cdot k} = \tilde{g}_{\nu \cdot (k-1)+1}$. Now, note that $|\mathbf{c}|$ defined in the proposition is also periodic with the fundamental frequency $\nu$ such that $n/\nu \in \mathbb{N}$. If so, then the same holds for $\mathbf{z}$ defined by (110) because adding a constant or rectifying a function does not change its periodicity properties. Combining this result with (114) validates the claim of (110).

To show (111), consider $\mathbf{W}_\varpi \mathbf{F}(\mathbf{m} \circ \mathbf{z})$. For every $r \in \mathcal{I}_n$, we have

$$
\begin{aligned}
(\mathbf{W}_\varpi \mathbf{F}(\mathbf{m} \circ \mathbf{z}))_r &= \frac{(\mathbf{W}_\varpi)_{rr}}{\sqrt{n}} \sum_{j=1}^n \left( m_j \cdot z_j \cdot e^{-\imath 2\pi (r-1)(j-1)/n} \right) \\
&= \frac{(\mathbf{W}_\varpi)_{rr}}{\sqrt{n}} \sum_{j=1}^n \left( m_j \cdot \left[ \sum_{l=1}^n \left( z_l \cdot \underbrace{\frac{1}{n} \sum_{k=1}^n e^{\imath 2\pi (k-1)(j-l)/n}}_{\delta_{j,l}} \right) \right] \cdot e^{-\imath 2\pi (r-1)(j-1)/n} \right) \\
&= \frac{(\mathbf{W}_\varpi)_{rr}}{\sqrt{n}} \sum_{k=1}^n \left( \frac{1}{\sqrt{n}} \sum_{l=1}^n \left( z_l \cdot e^{-\imath 2\pi (k-1)(l-1)/n} \right) \cdot \frac{1}{\sqrt{n}} \sum_{j=1}^n \left( m_j \cdot e^{-\imath 2\pi (j-1)(r-k)/n} \right) \right) \\
&= \frac{(\mathbf{W}_\varpi)_{rr}}{\sqrt{n}} \sum_{k=1}^n \left( \tilde{\tilde{z}}_k \cdot \tilde{\tilde{m}}_{r-k} \right) \\
&= \frac{\tilde{\tilde{z}}_1 \cdot (\mathbf{W}_\varpi)_{rr}}{\sqrt{n}} \cdot \tilde{\tilde{m}}_r + \underbrace{\frac{(\mathbf{W}_\varpi)_{rr}}{\sqrt{n}} \sum_{k=\nu+1}^{n-\nu+1} \left( \tilde{\tilde{z}}_k \cdot \tilde{\tilde{m}}_{r-k} \right)}_{=0 \;\Longleftarrow\; \omega \leqslant \varpi \leqslant \nu-\omega+1} \\
&= \langle \mathbf{z} \rangle \cdot \tilde{\tilde{m}}_r.
\end{aligned}
\tag{115}
$$

When writing the second equality above, we used the orthonormality of $\mathbf{F}^{-1}$. Combining (115) with the definition of $\mathbf{P}_{\mathcal{S}_\varpi}[\ldots]$ (see (10) in the main text), we obtain (111).

After establishing (110) and (111), consider the sequences $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ and $\mathbf{a}^{(0)}, \mathbf{a}^{(1)}, \ldots, \mathbf{a}^{(i)}, \ldots$. Note that

$$
\begin{aligned}
\mathbf{m}^{(0)} &= \mathbf{m} \circ |\mathbf{c}| \\
&= \mathbf{m} \circ \left( q^{(0)} + (|\mathbf{c}| - q^{(0)}) \circ \theta(|\mathbf{c}| - q^{(0)}) \right),
\end{aligned}
\tag{116}
$$

where $q^{(0)} = 0$. Hence, $\mathbf{m}^{(0)}$ can be expressed as an elementwise product of the true modulator $\mathbf{m}$ and a vector that satisfies (110). Let us now assume that, for some $i \geqslant 1$, $\mathbf{m}^{(i)}$ can be expressed as

$$
\mathbf{m}^{(i)} = \mathbf{m} \circ \left( q^{(i)} + (|\mathbf{c}| - q^{(i)}) \circ \theta(|\mathbf{c}| - q^{(i)}) \right).
\tag{117}
$$

Then, by (111),

$$
\mathbf{a}^{(i+1)} = \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(\mathbf{i})}] = q^{(i+1)} \cdot \mathbf{m},
\tag{118}
$$

where,

$$
q^{(i+1)} = \langle q^{(i)} + (|\mathbf{c}| - q^{(i)}) \circ \theta(|\mathbf{c}| - q^{(i)}) \rangle.
\tag{119}
$$

Further, by the definition of $\mathbf{P}_{\mathcal{S}_{\geqslant |\mathbf{S}|}}[\ldots]$ (see (9) in the main text),

$$
\begin{aligned}
\mathbf{m}^{(i+1)} &= \mathbf{P}_{\mathcal{S}_{\geqslant |\mathbf{S}|}}[\mathbf{a}^{(i+1)}] \\
&= \mathbf{m} \circ \left( q^{(i+1)} + (|\mathbf{c}| - q^{(i+1)}) \circ \theta(|\mathbf{c}| - q^{(i+1)}) \right),
\end{aligned}
\tag{120}
$$

i.e., $\mathbf{m}^{(i+1)}$ can also be expressed as the product of $\mathbf{m}$ and a vector that satisfies (110). Hence, we conclude by using mathematical induction that, for every $i \geqslant 1$,

$$\mathbf{a}^{(i)} = q^{(i)} \cdot \mathbf{m}, \tag{121}$$

and

$$q^{(i)} = \langle q^{(i-1)} + (|\mathbf{c}| - q^{(i-1)}) \circ \theta(|\mathbf{c}| - q^{(i-1)}) \rangle, \tag{122}$$

with $q_0 = 0$.

Next, observe that, by (122),

$$q^{(i)} - q^{(i-1)} = \langle (|\mathbf{c}| - q^{(i-1)}) \circ \theta(|\mathbf{c}| - q^{(i-1)}) \rangle$$
$$\geqslant 0, \tag{123}$$

i.e., the sequence $q^{(0)}, q^{(1)}, \ldots, q^{(i)}, \ldots$ is monotonically increasing. Moreover, it follows from (122) that, for every $i \geqslant 1$,

$$q^{(i-1)} \leqslant 1 \implies q^{(i)} = q^{(i-1)} + \langle (|\mathbf{c}| - q^{(i-1)}) \circ \theta(|\mathbf{c}| - q^{(i-1)}) \rangle$$
$$\leqslant q^{(i-1)} + (1 - q^{(i-1)}) \cdot \underbrace{\theta(1 - q^{(i-1)})}_{=1} = 1. \tag{124}$$

Taken together with $q^{(0)} = 0$, (124) implies that the sequence $q^{(0)}, q^{(1)}, \ldots, q^{(i)}, \ldots$ is bounded from above by 1. Hence, by the monotone convergence theorem (see *Proposition D.10*), $q^{(0)}, q^{(1)}, \ldots, q^{(i)}, \ldots$ converges to its least upper bound $\bar{q} \leqslant 1$. If we assume that $\bar{q} < 1$, then the convergence of $q^{(0)}, q^{(1)}, \ldots, q^{(i)}, \ldots$ to $\bar{q}$ implies that there exists an $N$ such that

$$\bar{q} - q^{(N)} \leqslant \langle (|\mathbf{c}| - \bar{q}) \circ \theta(|\mathbf{c}| - \bar{q}) \rangle / 2. \tag{125}$$

By (122),

$$q^{(N+1)} - q^{(N)} = \langle (|\mathbf{c}| - q^{(N)}) \circ \theta(|\mathbf{c}| - q^{(N)}) \rangle$$
$$\geqslant \langle (|\mathbf{c}| - \bar{q}) \circ \theta(|\mathbf{c}| - \bar{q}) \rangle$$
$$> \langle (|\mathbf{c}| - \bar{q}) \circ \theta(|\mathbf{c}| - \bar{q}) \rangle / 2, \tag{126}$$

which means that $q^{(N+1)} > \bar{q}$, i.e., the initial assumption that $\bar{q} < 1$ is incorrect. Therefore, $\bar{q} = 1$, and $q^{(0)}, q^{(1)}, \ldots, q^{(i)}, \ldots$ converges to 1, i.e.,

$$\forall \epsilon > 0 \; \exists N(\epsilon) : i > N(\epsilon) \implies |1 - q^{(N)}| < \epsilon. \tag{127}$$

Finally, note that

$$|1 - q^{(N)}| < \epsilon \implies \|\mathbf{m}\|_2 \cdot |1 - q^{(N)}| < \underbrace{\|\mathbf{m}\|_2 \cdot \epsilon}_{\epsilon'}$$
$$\implies \|\mathbf{m} - \mathbf{m} \cdot q^{(N)}\|_2 < \epsilon'$$
$$\implies \underbrace{\|\mathbf{m} - \mathbf{a}^{(N)}\|_2 < \epsilon'}_{\text{by (121)}}, \tag{128}$$

i.e., (127) implies that the sequence $\mathbf{a}^{(0)}, \mathbf{a}^{(1)}, \ldots, \mathbf{a}^{(i)}, \ldots$ converges to $\mathbf{m}$. In the light of (81) and (82) in the proof of *Proposition III.1*, this result allows concluding that $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ also converges to $\mathbf{m}$. ∎

### *AP-A algorithm*

**Proposition III.3.** *A sequence* $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ *formed by the AP-A algorithm for* $\epsilon_{tol} = 0$ *and* $N_{iter} \to +\infty$ *converges to some* $\mathbf{m}^{\dagger} \in \mathcal{S}_{\geqslant |\mathbf{s}|} \cap \mathcal{S}_{\varpi}$. *The convergence is monotonic, i.e.,* $\|\mathbf{m}^{(i+1)} - \mathbf{m}^{\dagger}\|_2 \leqslant \|\mathbf{m}^{(i)} - \mathbf{m}^{\dagger}\|_2$ *for* $i \geqslant 0$.

*Proof.* Analogously to the AP-B algorithm, if the sequence $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ terminates with some finite $i = N$, then we have $\mathbf{m}^{(N)} = \mathbf{m}^{\dagger} \in \mathcal{S}_{\geqslant |\mathbf{s}|} \cap \mathcal{S}_{\varpi}$. Therefore, we next consider the case when the sequence is infinite. The main idea behind the proof is to show that the inequalities (78) and (80) apply not only to the AP-B but also to the AP-A algorithm. When that is established, we can proceed along the path of the convergence proof of the AP-B scheme.

To this end, we first derive some auxiliary (in)equalities. Specifically, it follows from the definition of the operator $\mathbf{P}_{\mathcal{S}_{\varpi}}$ [see (74)] that, for all $\mathbf{z}, \mathbf{y} \in \mathbb{R}^n$,

$$
\begin{aligned}
\langle \mathbf{z}, \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{y}] \rangle &= \mathbf{z}^{\mathrm{T}} \mathbf{F}^{-1} \mathbf{W}_{\varpi} \mathbf{F} \mathbf{y} = \mathbf{z}^{\mathrm{T}} \mathbf{F}^{-1} \mathbf{W}_{\varpi} \mathbf{W}_{\varpi} \mathbf{F} \mathbf{y} = \mathbf{z}^{\mathrm{T}} \mathbf{F}^{-1} \mathbf{W}_{\varpi} \mathbf{F} \mathbf{F}^{-1} \mathbf{W}_{\varpi} \mathbf{F} \mathbf{y} \\
&= (\mathbf{z}^{\mathrm{T}} \mathbf{F}^{-1} \mathbf{W}_{\varpi} \mathbf{F})^{*} (\mathbf{F}^{-1} \mathbf{W}_{\varpi} \mathbf{F} \mathbf{y}) = (\mathbf{z}^{\mathrm{T}} \mathbf{F} \mathbf{W}_{\varpi} \mathbf{F}^{-1})(\mathbf{F}^{-1} \mathbf{W}_{\varpi} \mathbf{F} \mathbf{y}) \\
&= ((\mathbf{F} \mathbf{W}_{\varpi} \mathbf{F}^{-1})^{\mathrm{T}} \mathbf{z})^{\mathrm{T}} (\mathbf{F}^{-1} \mathbf{W}_{\varpi} \mathbf{F} \mathbf{y}) = (\mathbf{F}^{-1} \mathbf{W}_{\varpi} \mathbf{F} \mathbf{z})^{\mathrm{T}} (\mathbf{F}^{-1} \mathbf{W}_{\varpi} \mathbf{F} \mathbf{y}) \\
&= \langle \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{z}], \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{y}] \rangle .
\end{aligned}
\tag{129}
$$

Here, T and * mark, respectively, the transposition and complex conjugation. We used the following properties of matrices $\mathbf{W}_{\varpi}$ and $\mathbf{F}$ in (129): 1) $\mathbf{W}_{\varpi} \mathbf{W}_{\varpi} = \mathbf{W}_{\varpi}$; 2) $\mathbf{W}_{\varpi}^{*} = \mathbf{W}_{\varpi}^{\mathrm{T}} = \mathbf{W}_{\varpi}$; 3) $\mathbf{F}^{*} = \mathbf{F}^{-1}$; and 4) $\mathbf{F}^{\mathrm{T}} = \mathbf{F}$. The result of (129) can be rewritten as

$$
\langle \mathbf{z} - \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{z}], \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{y}] \rangle = 0.
\tag{130}
$$

(130) is a particular instance of a more general result that the difference between any $\mathbf{z} \in \mathbb{R}^n$ and its projection onto a linear subspace of $\mathbb{R}^n$ is perpendicular to any element of that subspace. Using (130), we obtain the following:

$$
\begin{aligned}
\|\mathbf{z}\|_2^2 = \langle \mathbf{z}, \mathbf{z} \rangle &= \langle \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{z}] + (\mathbf{z} - \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{z}]), \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{z}] + (\mathbf{z} - \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{z}]) \rangle \\
&= \|\mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{z}]\|_2^2 + \|\mathbf{z} - \mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{z}]\|_2^2 \\
&\geqslant \|\mathbf{P}_{\mathcal{S}_{\varpi}}[\mathbf{z}]\|_2^2 .
\end{aligned}
\tag{131}
$$

Applying (131) to the *line 6* of the AP-A algorithm $\left( \lambda = \|\mathbf{m}^{(i-1)} - \mathbf{a}^{(i-1)}\|_2^2 / \|\mathbf{b}^{(i)}\|_2^2 \right)$, we get

$$
\lambda \geqslant 1,
\tag{132}
$$

with the equality holding if and only if $(\mathbf{m}^{(i-1)} - \mathbf{a}^{(i-1)})$ is mapped by $\mathbf{P}_{\mathcal{S}_\varpi}$ to itself, i.e., $\mathbf{m}^{(i-1)} \in \mathcal{S}_\varpi$. However, this would mean that the convergence was reached at iteration $i-1$. Finally, we note that *line 7* of the AP-A algorithm $\left(\mathbf{a}^{(i)} = \mathbf{a}^{(i-1)} + \lambda \cdot \mathbf{b}^{(i)}\right)$ implies

$$(\mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}] - \mathbf{a}^{(i+1)}) = (\lambda - 1) \cdot (\mathbf{a}^{(i)} - \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}]) \tag{133}$$

and

$$
\begin{aligned}
\langle \mathbf{m}^{(i)} - \mathbf{a}^{(i+1)}, \mathbf{m}^{(i)} - \mathbf{a}^{(i)} \rangle &= \langle \mathbf{m}^{(i)} - \mathbf{a}^{(i)}, \mathbf{m}^{(i)} - \mathbf{a}^{(i)} \rangle - \lambda \cdot \langle \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}] - \mathbf{a}^{(i+1)}, \mathbf{m}^{(i)} - \mathbf{a}^{(i)} \rangle \\
&= \|\mathbf{m}^{(i)} - \mathbf{a}^{(i)}\|_2^2 - \frac{\|\mathbf{m}^{(i)} - \mathbf{a}^{(i)}\|_2^2}{\|\mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}] - \mathbf{a}^{(i)}\|_2^2} \cdot \|\mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}] - \mathbf{a}^{(i)}\|_2^2 \\
&= 0.
\end{aligned}
\tag{134}
$$

In (134) we applied (129) to the term $\langle \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}] - \mathbf{a}^{(i+1)}, \mathbf{m}^{(i)} - \mathbf{a}^{(i)} \rangle$.

We are now ready to prove that (78) and (80) hold for the AP-A algorithm. In particular, we have

$$
\begin{aligned}
\|\mathbf{x} - \mathbf{m}^{(i)}\|_2^2 &= \|\mathbf{x} - \mathbf{a}^{(i+1)} + \mathbf{a}^{(i+1)} - \mathbf{m}^{(i)}\|_2^2 \\
&= \|\mathbf{x} - \mathbf{a}^{(i+1)}\|_2^2 + \|\mathbf{a}^{(i+1)} - \mathbf{m}^{(i)}\|_2^2 + 2 \cdot \langle \mathbf{m}^{(i)} - \mathbf{a}^{(i+1)}, \mathbf{a}^{(i+1)} - \mathbf{x} \rangle.
\end{aligned}
\tag{135}
$$

The first two terms in the second line of (135) are nonnegative by the definition of the norm. The last term is nonnegative too (note that $\mathbf{x} \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_\varpi$):

$$
\begin{aligned}
\langle \mathbf{m}^{(i)} - \mathbf{a}^{(i+1)}, \mathbf{a}^{(i+1)} - \mathbf{x} \rangle &= \langle \mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}] + \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}] - \mathbf{a}^{(i+1)}, \mathbf{a}^{(i+1)} - \mathbf{x} \rangle \\
&= \underbrace{\langle \mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}], \mathbf{a}^{(i+1)} - \mathbf{x} \rangle}_{=0 \text{ by (130)}} + \langle \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}] - \mathbf{a}^{(i+1)}, \mathbf{a}^{(i+1)} - \mathbf{x} \rangle \\
&= \underbrace{(\lambda - 1) \cdot \langle \mathbf{a}^{(i)} - \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}], \mathbf{a}^{(i+1)} - \mathbf{x} \rangle}_{\text{by (133)}} \\
&= (\lambda - 1) \cdot \langle \mathbf{a}^{(i)} - \mathbf{m}^{(i)} + \mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}], \mathbf{a}^{(i+1)} - \mathbf{x} \rangle \\
&= (\lambda - 1) \cdot \underbrace{\langle \mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}], \mathbf{a}^{(i+1)} - \mathbf{x} \rangle}_{=0 \text{ by (130)}} + (\lambda - 1) \cdot \langle \mathbf{a}^{(i)} - \mathbf{m}^{(i)}, \mathbf{a}^{(i+1)} - \mathbf{x} \rangle \\
&= (\lambda - 1) \cdot \langle \mathbf{a}^{(i)} - \mathbf{m}^{(i)}, \mathbf{a}^{(i+1)} - \mathbf{m}^{(i)} + \mathbf{m}^{(i)} - \mathbf{x} \rangle \\
&= (\lambda - 1) \cdot \underbrace{\langle \mathbf{a}^{(i)} - \mathbf{m}^{(i)}, \mathbf{a}^{(i+1)} - \mathbf{m}^{(i)} \rangle}_{=0 \text{ by (134)}} + \underbrace{(\lambda - 1)}_{\geqslant 0 \text{ by (132)}} \cdot \underbrace{\langle \mathbf{a}^{(i)} - \mathbf{m}^{(i)}, \mathbf{m}^{(i)} - \mathbf{x} \rangle}_{\geqslant 0 \text{ by (65)}} \\
&\geqslant 0.
\end{aligned}
\tag{136}
$$

Therefore,

$$\|\mathbf{x} - \mathbf{m}^{(i)}\|_2^2 \geqslant \|\mathbf{x} - \mathbf{a}^{(i+1)}\|_2^2 + \|\mathbf{a}^{(i+1)} - \mathbf{m}^{(i)}\|_2^2. \tag{137}$$

The derivation of the inequality (80) for the AP-A algorithm is equivalent to that for the AP-B:

$$
\begin{aligned}
\|\mathbf{x} - \mathbf{a}^{(i+1)}\|_2^2 &= \|\mathbf{x} - \mathbf{m}^{(i+1)} + \mathbf{m}^{(i+1)} - \mathbf{a}^{(i+1)}\|_2^2 \\
&= \|\mathbf{x} - \mathbf{m}^{(i+1)}\|_2^2 + \|\mathbf{m}^{(i+1)} - \mathbf{a}^{(i+1)}\|_2^2 + 2 \cdot \underbrace{\left\langle \mathbf{a}^{(i+1)} - \mathbf{m}^{(i+1)}, \mathbf{m}^{(i+1)} - \mathbf{x} \right\rangle}_{\geqslant 0 \text{ by (65)}}.
\end{aligned}
\tag{138}
$$

Thus,

$$
\|\mathbf{x} - \mathbf{a}^{(i+1)}\|_2^2 \geqslant \|\mathbf{x} - \mathbf{m}^{(i+1)}\|_2^2 + \|\mathbf{m}^{(i+1)} - \mathbf{a}^{(i+1)}\|_2^2.
\tag{139}
$$

The rest of the proof follows step by step the proof of *Proposition III.1*. Indeed, starting with the inequalities (78) and (80), to which (137) and (139) are equivalent, the proof of *Proposition III.1* proceeds based on them and the closedness and convexity of the sets $\mathcal{S}_{\geqslant |\mathbf{s}|}$ and $\mathcal{S}_\varpi$ entirely. The monotonicity of the convergence follows from an equivalent of (90), which is derived as a part of the proof of the convergence itself.

∎

*Remark.* 1) Concerning the set $\mathcal{S}_{\geqslant |\mathbf{s}|}$, the proof above relies entirely on its closedness and convexity. Thus, the AP-A algorithm is still valid under these more general assumptions. 2) The expressions (129)$-$(134) apply to projection operators onto any linear space. Hence, the AP-A algorithm still works if $\mathcal{S}_\varpi$ is replaced by an arbitrary linear space.

**Proposition III.4.** *Consider* $\mathbf{m} \in \mathcal{M}_\omega$ *and* $\mathbf{c} \in \mathcal{C}_d$ *with* $|c_j| = \sum_{k=1}^{n/\nu} (\tilde{c}_{\nu \cdot k} \cdot e^{i 2 \pi \nu (k-1)(j-1)/n})$, *where* $\tilde{c}_{\nu \cdot k} \in \mathbb{C}$ *and* $n/\nu \in \mathbb{N}$. *If* $\varpi \geqslant \omega$ *and* $\nu \geqslant \varpi + \omega - 1$, *then a sequence* $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ *formed by the AP-A algorithm for* $\epsilon_{tol} = 0$ *and* $N_{iter} \to +\infty$ *converges to* $\mathbf{m}$.

*Proof.* The proof of this proposition goes along the lines of that of *Proposition III.2*. First, by using (110) and (111), we derive the result analogous to (121) and (122). Then, we show the monotonic convergence of $q^{(0)}, q^{(1)}, \ldots, q^{(i)}, \ldots$ to 1, which assures the convergence of $\mathbf{a}^{(0)}, \mathbf{a}^{(1)}, \ldots, \mathbf{a}^{(i)}, \ldots$ and $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ to $\mathbf{m}$.

To derive the result analogous to (121) and (122), assume that, for some $i \geqslant 0$ and $q^{(i)} \in \mathbb{R}$,

$$
\begin{aligned}
\mathbf{a}^{(i)} &= q^{(i)} \cdot \mathbf{m}, \\
\mathbf{m}^{(i)} &= \mathbf{m} \circ \left( q^{(i)} + (|\mathbf{c}| - q^{(i)}) \circ \theta(|\mathbf{c}| - q^{(i)}) \right).
\end{aligned}
\tag{140}
$$

First, note that (140) applies when $i = 0$ with $q^{(0)} = 0$. Next, by the definition of the AP-A algorithm, we have

$$
\begin{aligned}
\mathbf{b}^{(i+1)} &= \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)} - \mathbf{a}^{(i)}] \\
&= \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}] - \mathbf{a}^{(i)} \\
&= \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m} \circ (q^{(i)} + (|\mathbf{c}| - q^{(i)}) \circ \theta(|\mathbf{c}| - q^{(i)}))] - q^{(i)} \cdot \mathbf{m} \\
&= \underbrace{\left\langle (|\mathbf{c}| - q^{(i)}) \circ \theta(|\mathbf{c}| - q^{(i)}) \right\rangle \cdot \mathbf{m}}_{\text{by (110) and (111)}},
\end{aligned}
\tag{141}
$$

$$\lambda = \frac{\|\mathbf{m}^{(i)} - \mathbf{a}^{(i)}\|_2}{\|\mathbf{b}^{(i+1)}\|_2}$$

$$= \frac{\|(|\mathbf{c}| - q^{(i)}) \circ \theta(|\mathbf{c}| - q^{(i)}) \circ \mathbf{m}\|_2}{\langle (|\mathbf{c}| - q^{(i)}) \circ \theta(|\mathbf{c}| - q^{(i)}) \rangle \cdot \|\mathbf{m}\|_2}, \tag{142}$$

$$\mathbf{a}^{(i+1)} = \mathbf{a}^{(i)} + \lambda \cdot \mathbf{b}^{(i+1)}$$

$$= \underbrace{\left( q^{(i)} + \frac{\|(|\mathbf{c}| - q^{(i)}) \circ \theta(|\mathbf{c}| - q^{(i)}) \circ \mathbf{m}\|_2}{\|\mathbf{m}\|_2} \right)}_{q^{(i+1)}} \cdot \mathbf{m}, \tag{143}$$

and

$$\mathbf{m}^{(i+1)} = \mathbf{P}_{\mathcal{S}_{\geqslant |\mathbf{S}|}}[\mathbf{a}^{(i+1)}],$$

$$= \mathbf{m} \circ \left( q^{(i+1)} + (|\mathbf{c}| - q^{(i+1)}) \circ \theta(|\mathbf{c}| - q^{(i+1)}) \right). \tag{144}$$

Applying the principle of mathematical induction to the above results, we conclude that, for any $i \geqslant 0$,

$$\mathbf{a}^{(i)} = q^{(i)} \cdot \mathbf{m}, \tag{145}$$

$$q^{(i)} = \left( q^{(i-1)} + \frac{\|(|\mathbf{c}| - q^{(i-1)}) \circ \theta(|\mathbf{c}| - q^{(i-1)}) \circ \mathbf{m}\|_2}{\|\mathbf{m}\|_2} \right), \tag{146}$$

with $q^{(0)} = 0$.

By (146),

$$q^{(i)} - q^{(i-1)} = \frac{\|(|\mathbf{c}| - q^{(i-1)}) \circ \theta(|\mathbf{c}| - q^{(i-1)}) \circ \mathbf{m}\|_2}{\|\mathbf{m}\|_2} \tag{147}$$

$$\geqslant 0,$$

i.e., the sequence $q^{(0)}, q^{(1)}, \ldots$ is monotonically increasing. Moreover, it follows from (146) that, for every $i \geqslant 1$,

$$q^{(i-1)} \leqslant 1 \implies q^{(i)} = \left( q^{(i-1)} + \frac{\|(|\mathbf{c}| - q^{(i-1)}) \circ \theta(|\mathbf{c}| - q^{(i-1)}) \circ \mathbf{m}\|_2}{\|\mathbf{m}\|_2} \right)$$

$$\leqslant \left( q^{(i-1)} + \frac{\|(1 - q^{(i-1)}) \cdot \mathbf{m}\|_2}{\|\mathbf{m}\|_2} \right) = 1, \tag{148}$$

Taken together with $q^{(0)} = 0$, (148) implies that the sequence $q^{(0)}, q^{(1)}, \ldots, q^{(i)}, \ldots$ is bounded from above by 1. Hence, by the monotone convergence theorem (see *Proposition D.10*), $q^{(0)}, q^{(1)}, \ldots, q^{(i)}, \ldots$ converges to its least upper bound $\bar{q} \leqslant 1$. If we assume that $\bar{q} < 1$, then the convergence of $q^{(0)}, q^{(1)}, \ldots, q^{(i)}, \ldots$ to $\bar{q}$ implies that there exists an $N$ such that

$$\bar{q} - q^{(N)} \leqslant \frac{1}{2} \cdot \frac{\|(|\mathbf{c}| - \bar{q}) \circ \theta(|\mathbf{c}| - \bar{q}) \circ \mathbf{m}\|_2}{\|\mathbf{m}\|_2}. \tag{149}$$

By (146),

$$q^{(N+1)} - q^{(N)} = \frac{\|(|\mathbf{c}| - q^{(N)}) \circ \theta(|\mathbf{c}| - q^{(N)}) \circ \mathbf{m}\|_2}{\|\mathbf{m}\|_2}$$

$$\geqslant \frac{\|(|\mathbf{c}| - \bar{q}) \circ \theta(|\mathbf{c}| - \bar{q}) \circ \mathbf{m}\|_2}{\|\mathbf{m}\|_2} \tag{150}$$

$$> \frac{1}{2} \cdot \frac{\|(|\mathbf{c}| - \bar{q}) \circ \theta(|\mathbf{c}| - \bar{q}) \circ \mathbf{m}\|_2}{\|\mathbf{m}\|_2},$$

which means that $q^{(N+1)} > \bar{q}$, i.e., the initial assumption that $\bar{q} < 1$ is incorrect. Therefore, $\bar{q} = 1$, and $q^{(0)}, q^{(1)}, \ldots, q^{(i)}, \ldots$ converges to 1. This, as shown in the last paragraph of the proof of *Proposition III.2*, implies that $\mathbf{a}^{(0)}, \mathbf{a}^{(1)}, \ldots, \mathbf{a}^{(i)}, \ldots$ and $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ converge to $\mathbf{m}$.                                 ∎

### *AP-P algorithm*

**Proposition III.5.** *A sequence* $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ *formed by the AP-P algorithm for* $\epsilon_{tol} = 0$ *and* $N_{iter} \to +\infty$ *converges to a unique* $\mathbf{m}^\dagger \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$ *such that* $\|\mathbf{m}^\dagger\|_2 \leqslant \|\mathbf{x}\|_2$ *for every* $\mathbf{x} \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$. *The convergence is monotonic, i.e.,* $\|\mathbf{m}^{(i+1)} - \mathbf{m}^\dagger\|_2 \leqslant \|\mathbf{m}^{(i)} - \mathbf{m}^\dagger\|_2$ *for* $i \geqslant 0$.

*Proof.* The proof follows as a corollary of a more general theorem proved for a finite number of closed convex sets in a Hilbert space by Boyle and Dykstra [6, *Theorem 2*]. Specifically, for $\mathbf{m}^{(0)} = \mathbf{0}$ and particular constraint sets $\mathcal{S}_{\geqslant|\mathbf{s}|}$ and $\mathcal{S}_{\varpi}$, the sequence $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ formed by the algorithm formulated there (the so-called Dykstra's algorithm) converges to a unique $\mathbf{m}^\dagger \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$ such that $\|\mathbf{m}^\dagger\|_2 \leqslant \|\mathbf{x}\|_2$ for any $\mathbf{x} \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$. For our purposes, it is thus enough to show that the sequence $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ generated by the AP-P algorithm converges to the same $\mathbf{m}^\dagger$.

First, we observe that Dykstra's algorithm formally turns into the AP-P (except the difference in the initial conditions) if we set r = 2, denote $g_{i,1} \equiv \mathbf{a}^{(i-1)}$, $g_{i,2} \equiv \mathbf{m}^{(i-1)}$, $I_{i,2} \equiv \mathbf{c}^{(i-1)}$, and assign $I_{i,1} \equiv 0$ there (see *Theorem 2* in [6]). Note that, originally, $I_{i,1} = g_{i,1} - (g_{i,2} - I_{i-1,1})$ and $g_{i,1} = \mathbf{P}_{\mathcal{S}_{\varpi}}[g_{i-1,2} - I_{i-1,1}]$. However, due to the linearity of $\mathbf{P}_{\mathcal{S}_{\varpi}}$ [see (74)], we have

$$
\begin{aligned}
g_{i,1} &= \mathbf{P}_{\mathcal{S}_{\varpi}}[g_{i-1,2}] + \mathbf{P}_{\mathcal{S}_{\varpi}}[I_{i-1,1}] \\
&= \mathbf{P}_{\mathcal{S}_{\varpi}}[g_{i-1,2}] + \underbrace{\mathbf{P}_{\mathcal{S}_{\varpi}}[g_{i-1,1}]}_{=g_{i-1,1}} - \underbrace{\mathbf{P}_{\mathcal{S}_{\varpi}}[g_{i-1,2}]}_{=g_{i-1,1}} + \mathbf{P}_{\mathcal{S}_{\varpi}}[I_{i-2,1}] \\
&= \mathbf{P}_{\mathcal{S}_{\varpi}}[g_{i-1,2}] + \mathbf{P}_{\mathcal{S}_{\varpi}}[I_{i-2,1}].
\end{aligned}
\tag{151}
$$

Applying (151) iteratively, we obtain $g_{i,1} = \mathbf{P}_{\mathcal{S}_{\varpi}}[g_{i-1,2}] + \mathbf{P}_{\mathcal{S}_{\varpi}}[I_{0,1}] = \mathbf{P}_{\mathcal{S}_{\varpi}}[g_{i-1,2}]$ because $I_{0,1} = 0$ by construction. Therefore, given the constraint sets $\mathcal{S}_{\geqslant|\mathbf{s}|}$ and $\mathcal{S}_{\varpi}$, $I_{i,1}$ can indeed be canceled from the Dykstra's algorithm by setting it to 0 without any consequences to its convergence properties.

Further, assuming $\mathbf{m}^{(0)} = \mathbf{c}^{(0)} = \mathbf{0}$, we can verify that Dykstra's algorithm after the first iteration coincides with the initialized AP-P algorithm, i.e., the latter is shifted forward by one iteration with respect to the former. Hence, the sequence $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(i)}, \ldots$ generated by the AP-P converges to the same $\mathbf{m}^\dagger \in \mathcal{S}_{\geqslant|\mathbf{s}|} \cap \mathcal{S}_{\varpi}$ as the analogous sequence formed by the Dykstra's algorithm initialized with $\mathbf{m}^{(0)} = \mathbf{0}$.

The monotonicity of the convergence of Dykstra's, and thus the AP-P, algorithms follows from the fact that all terms on the right-hand side of [6, (3.2)] are nonnegative and that each subsequent iteration only adds additional terms without discarding the old ones.                                 ∎

## G. Lower Bound on the Convergence Error

The proof of the statement about the lower bound on the convergence error made in Sections III-A and III-B of the main text are presented here.

**Proposition G.1.** *In the case of the AP-B and AP-A algorithms, the convergence error* $\|\mathbf{m}^{(i-1)} - \mathbf{m}^\dagger\|_2/\sqrt{n}$ *is bounded from below by the infeasibility error* $\epsilon^{(i)}$ *for any* $i \geqslant 1$.

*Proof.* According to (79), which applies to both the AP-B and AP-A algorithms,

$$\|\mathbf{m}^{(i-1)} - \mathbf{m}^\dagger\|_2/\sqrt{n} \geqslant \|\mathbf{a}^{(i)} - \mathbf{m}^\dagger\|_2/\sqrt{n} \qquad \forall i \geqslant 1. \tag{152}$$

Moreover, by the *Definition D.13* of the projection operator and the fact that $\mathbf{m}^\dagger \in \mathcal{S}_{\geqslant|\mathbf{s}|}$, we have

$$\|\mathbf{a}^{(i)} - \mathbf{m}^\dagger\|_2/\sqrt{n} \geqslant \|\mathbf{a}^{(i)} - \mathbf{P}_{\mathcal{S}_{\geqslant|\mathbf{s}|}}[\mathbf{a}^{(i)}]\|_2/\sqrt{n} = \|\mathbf{a}^{(i)} - \mathbf{m}^{(i)}\|_2/\sqrt{n} = \epsilon^{(i)} \qquad \forall i \geqslant 1. \tag{153}$$

Thus, $\|\mathbf{m}^{(i-1)} - \mathbf{m}^\dagger\|_2/\sqrt{n} \geqslant \epsilon^{(i)}$ for all $i \geqslant 1$. ∎

*Remark.* Using the *Definition D.13* of the projection operator and the fact that $\mathbf{m}^\dagger \in \mathcal{S}_\varpi$, we similarly conclude that

$$\|\mathbf{m}^{(i)} - \mathbf{m}^\dagger\|_2/\sqrt{n} \geqslant \|\mathbf{m}^{(i)} - \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{m}^{(i)}]\|_2/\sqrt{n} \qquad \forall i \geqslant 0. \tag{154}$$

(154) also applies in the context of the AP-P algorithm.

**SIGNAL CLASSES**

### H. TYPES OF WIDEBAND CARRIERS FOUND IN PRACTICE

Wideband carriers found in practice fall into three main classes: *1)* **(quasi-)harmonic**; *2)* **(quasi-)random**; and *3)* **spike-train**. Below, we provide examples of real-world signals featuring these carrier types and applications of their demodulation.

The need for demodulating signals formed of harmonic carriers is well recognized in acoustic imaging [7], [8]. There, sinusoidal wavepackets are used as probing signals. However, in many situations, the interaction between sound and matter is nonlinear. This makes the returning waves, whose time-dependent amplitude carries information about the imaged object, harmonic. The possibility of efficient and accurate demodulation of signals of this type would also allow generalizing the probing wave packets themselves from sinusoidal to harmonic.

A representative example of quasi-random carriers manifests in surface electromyography. In particular, an electrical signal detected at the skin surface during the skeletal muscle activity is an amplitude-modulated colored noise resulting from low-pass filtering of electric pulse trains generated by a large set of conditionally independent muscle fibers [9, Ch. 5]. The amplitude component of the recorded electrical signal carries information about the force pattern generated by the muscle being studied [10], [11] (see Fig. 11 for an example).
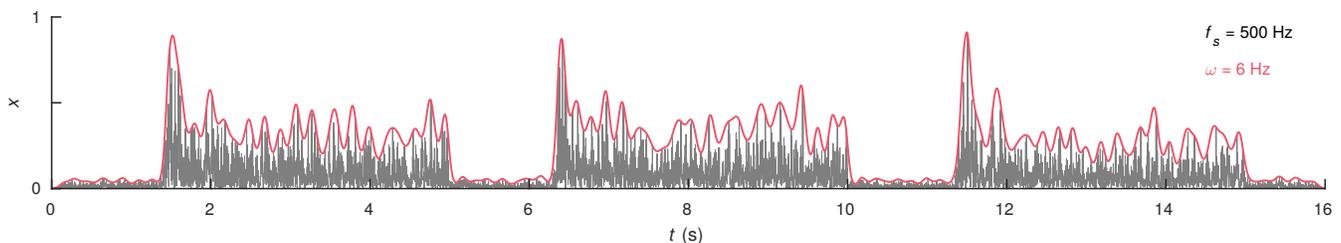


Fig. 11.  Demodulation of an electromyogram. The figure shows a 16 s long absolute-value electromyogram signal corresponding to the response of a flexor carpi ulnaris muscle (forearm) to three consecutive grasps of a spherical object (grey) and its modulator inferred by using the AP-A algorithm (red). The original recording was taken from [12], [13].

Probably the most elaborate applications of amplitude demodulation in the context of wideband signals are found in human speech processing. Speech signals are built of temporarily structured segments of quasi-harmonic and quasi-random carriers [14] that are amplitude-modulated at different timescales [15], [16]. Amplitude demodulation of broad-band, as well as sub-band speech demodulation, is widely exploited in automatic speech recognition [17], [18], [19], hearing restoration [20], [21] tasks, and fundamental studies of the neural mechanisms of auditory information processing in the brain [22], [23], [24], [25]. In all these cases, the modulators and carriers convey the information about specific aspects of speech, e.g., semantic meaning, associated emotion, or speaker identity, that need to be extracted. Demodulation of signals with mixed quasi-harmonic and quasi-random carrier types is also known in diagnostic phonocardiography [26], [27]. There, the extracted modulators of heart sounds are used for the detection of events of normal or abnormal functioning of this organ.

One of the most popular applications using demodulation of signals built of the spike-train type carriers is the Pulse-Code Modulation (PCM) protocol [28]. There, pulses are regularly spaced with specified locations and have a known constant amplitude and vanishing width.[17] More complex quasi-regular or stochastic pulse sequences of finite width manifest in electric activities of the cardiac muscles and neurons [9, Ch. 6], [29, Ch. 1]. Physiological and diagnostic information contained in these responses is typically associated with the microstructure and timing of the pulses. However, these recordings often come contaminated by artifacts or other physiological signals that slowly modulate the baseline and amplitude of the pulses [30], [9, Ch. 7.1]. Hence, to separate the fast cardiac or neural activity (the carrier) from other slowly changing physiological processes or artifacts (the modulator), the raw signal must be demodulated [30] (see Fig. 12 for an example). Finally, as we discuss in Sections IX.A and IX.B of the main text, carriers of the spike-train type also effectively manifest in applications when modulator recovery from nonuniformly or sparsely sampled signals of any origin is needed.
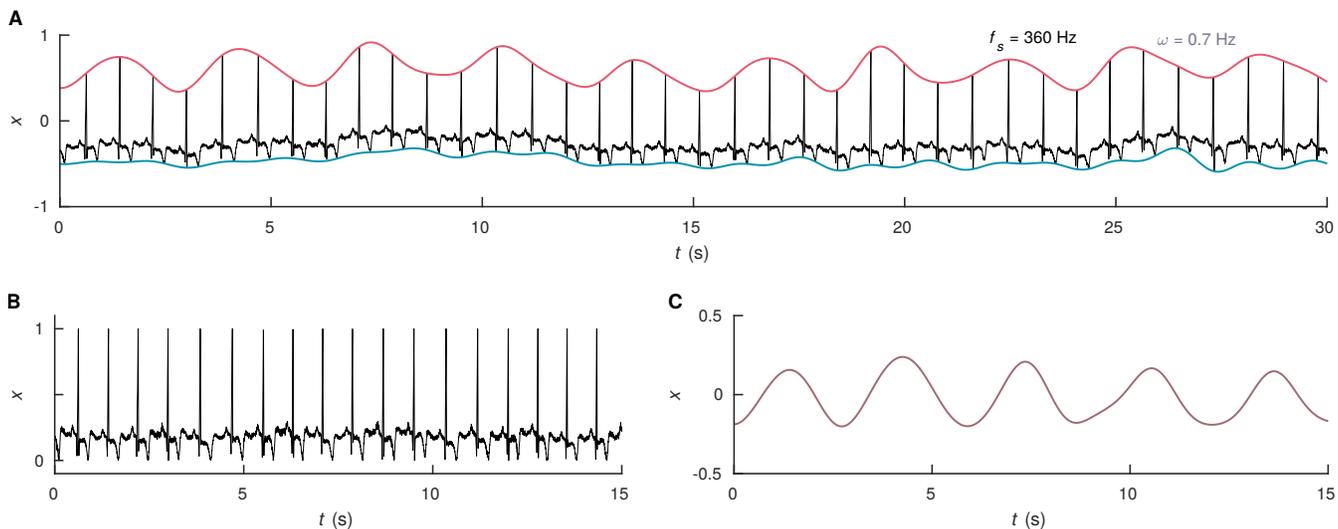


Fig. 12.   Demodulation of an electrocardiogram. **A**: A 30 s long fragment of an electrocardiogram recording (black) and its upper (red) and lower (blue) envelopes inferred by using the AP-A algorithm.[19] **B**: The first 15 s of the recovered carrier, i.e., the desired electrocardiogram with canceled artifacts. **C**: The first 15 s of the respiratory curve estimated from the upper and lower envelopes of the original signal in panel A. The original recording was taken from [31], [32].

## I. SYNTHETIC TEST SIGNALS

The mathematical models and numerical sampling algorithms of synthetic modulators and carriers examined in this work are described next.

---

[17]These properties enable the pulse-coded signals to be demodulated by a simple low-pass filtering procedure [28].

[19]Here, we define the upper envelope of a signal $\mathbf{s}$ as $\mathbf{m} + \min[\mathbf{s}] \cdot \mathbf{1}$, where $\mathbf{m}$ is the modulator of $\mathbf{s} - \min[\mathbf{s}] \cdot \mathbf{1}$. Accordingly, the lower envelope of $\mathbf{s}$ is defined as the negative upper envelope of $-\mathbf{s}$.

*Modulators (recovery tests)*

For signal recovery tests discussed in Section C, modulators were generated by uniformly sampling from $\mathcal{M}_\omega$. Without loss of generality, we assumed the subset of $\mathcal{M}_\omega$ whose elements' norms are fixed to $\sqrt{n}$. The sampling was achieved by using a specially adapted version of the Metropolis-Hastings algorithm (see [33, p. 411] for an introduction to this method). In the concrete, starting with some $\mathbf{m}^{(0)} \in \mathcal{M}_\omega$, a sequence $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots$ defined by

$$\mathbf{m}^{(i)} = \begin{cases} \sqrt{n} \cdot \mathbf{r}^{(i)}/\|\mathbf{r}^{(i)}\|_2, & \text{if } r_j^{(i)} \geqslant 0 \quad \forall j \in \mathcal{I}_n \\ \mathbf{m}^{(i-1)}, & \text{otherwise} \end{cases} \tag{155}$$

for $i \geqslant 1$ was generated. Here,

$$\mathbf{r}^{(i)} = \mathbf{m}^{(i-1)} + \mathbf{F}^{-1}\mathbf{g}^{(i)}, \tag{156}$$

and $\mathbf{g}^{(i)} \in \mathbb{C}^n$ such that

$$\begin{aligned}
&\text{Re}\left(g_1^{(i)}\right) \sim \mathcal{N}(0, \sigma), \quad \text{Im}\left(g_1^{(i)}\right) = 0, \\
&\text{Re}\left(g_j^{(i)}\right) \sim \mathcal{N}(0, \sigma), \quad \text{Im}\left(g_j^{(i)}\right) \sim \mathcal{N}(0, \sigma), \quad 2 \leqslant j \leqslant \omega, \\
&g_j^{(i)} = \left(g_{n+2-j}^{(i)}\right)^*, \quad n+2-\omega \leqslant j \leqslant n, \\
&g_j^{(i)} = 0, \quad \omega+1 \leqslant j \leqslant n+1-\omega,
\end{aligned} \tag{157}$$

and $g_j \perp\!\!\!\perp g_k$ for $j \neq k$. We adjusted the parameter $\sigma$ in (157) to achieve the acceptance rate of the sampling scheme (155) between 0.4 and 0.6. The initial $\mathbf{m}^{(0)}$ was taken as

$$\mathbf{m}^{(0)} = \sqrt{n} \cdot (\mathbf{g}^{(0)} - \min[\mathbf{g}^{(0)}] \cdot \mathbf{1})/\|\mathbf{g}^{(0)} - \min[\mathbf{g}^{(0)}] \cdot \mathbf{1}\|_2. \tag{158}$$

The iterative process generating the sequence $\mathbf{m}^{(0)}, \mathbf{m}^{(1)}, \ldots$ was terminated upon the first instance of adherence to the following equilibration criterion of the underlying Markov chain:

$$\frac{1}{2\omega - 2} \cdot \sum_{j=2}^{\omega} \left| \frac{1}{i} \cdot \sum_{l=1}^{i} (\mathbf{F}\mathbf{m}^{(l)})_j \right|^2 < 0.01 \cdot \frac{n}{2\omega - 1}. \tag{159}$$

The corresponding $\mathbf{m}^{(i)}$ was then chosen as a sample point from $\mathcal{M}_\omega$.

*Modulators (performance tests)*

For the performance, convergence, and robustness tests of the demodulation algorithm discussed in Sections IV – VI of the main text, two types of modulators were considered:

- *Nonstationary Gaussian* modulators were produced by transforming a delta-correlated Gaussian process with a time-dependent low-pass filter. Specifically, $\mathbf{m}$ was calculated as

$$\mathbf{m} = \boldsymbol{w} \circ (\mathbf{m}' - \min[\mathbf{m}'] \cdot \mathbf{1})/\max[\mathbf{m}' - \min[\mathbf{m}'] \cdot \mathbf{1}], \tag{160}$$

with $w$ being a window "function" (see (169)), and $\mathbf{m}'$ defined by

$$m'_i = \sum_{j=1}^{n} \left( g_j^{(0)} \cdot \left( \mathbf{F}^{-1} \mathbf{h}^{(i)} \right)_{i-j} \right), \quad i \in \mathcal{I}_n. \tag{161}$$

In (161), for every $i \in \mathcal{I}_n$,

$$
\begin{aligned}
&\mathrm{Re}\left( h_1^{(i)} \right) \sim \left( \mathbf{P}_{\mathcal{S}_\omega}[\mathbf{g}^{(1)}] \right)_i, \quad \mathrm{Im}\left( h_1^{(i)} \right) = 0, \\
&\mathrm{Re}\left( h_j^{(i)} \right) = \left( \mathbf{P}_{\mathcal{S}_\omega}[\mathbf{g}^{(2j-2)}] \right)_i, \quad \mathrm{Im}\left( h_j^{(i)} \right) = \left( \mathbf{P}_{\mathcal{S}_\omega}[\mathbf{g}^{(2j-1)}] \right)_i, \quad 2 \leqslant j \leqslant \omega, \\
&h_j^{(i)} = \left( h_{n+2-j}^{(i)} \right)^*, \quad n+2-\omega \leqslant j \leqslant n, \\
&h_j^{(i)} = 0, \quad \omega + 1 \leqslant j \leqslant n+1-\omega,
\end{aligned}
\tag{162}
$$

with $\mathbf{g}^{(j)}$, $j \in \mathcal{I}_{2\omega-1}$, being independent samples from the standard $n$-dimensional Gaussian distribution, i.e., $g_l^{(j)} \sim \mathcal{N}(0,1)$ for $l \in \mathcal{I}_n$, and $g_l^{(j)} \perp\!\!\!\perp g_k^{(j)}$ for $l \neq k$.

- *Maximally-uniformly distributed* modulators were created by using the NORTA algorithm [34]. Specifically, the modulators were calculated as

$$\mathbf{m} = \boldsymbol{w} \circ (\mathbf{m}' - \min[\mathbf{m}'] \cdot \mathbf{1}) / \max[\mathbf{m}' - \min[\mathbf{m}'] \cdot \mathbf{1}], \tag{163}$$

with $w$ being the window "function" (see (169)), and $\mathbf{m}'$ given by

$$\mathbf{m}' = \mathbf{P}_{\mathcal{S}_\omega}[\Phi(\mathbf{F}^{-1}\mathbf{r})], \tag{164}$$

where $\Phi(\dots)$ is the cumulative distribution function of the standard Gaussian random variable and $\mathbf{r}$ is given by

$$
r_j = \begin{cases}
\sqrt{p_1^\star} \cdot g_1, & j = 1 \\
\sqrt{p_j^\star/2} \cdot (g_{2j-2} + \imath \cdot g_{2j-1}), & 2 \leqslant j \leqslant \lfloor (n+1)/2 \rfloor \\
\sqrt{p_{(n+2)/2}^\star} \cdot g_n, & j = (n+2)/2 \\
(r_{n+2-j})^*, & (n+2)/2 < j \leqslant n
\end{cases}. \tag{165}
$$

In (165), $\mathbf{g}$ is the standard $n$-dimensional Gaussian random vector, and

$$\mathbf{p}^\star = \sqrt{n} \cdot (\mathbf{F}\mathbf{c}^\star) \circ \theta(\mathbf{F}\mathbf{c}^\star), \tag{166}$$

with $c_1^\star = 1$, and the remaining components of $\mathbf{c}^\star$ implicitly defined by

$$(\mathbf{F}^{-1}\mathbf{p}/\sqrt{n})_j = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \left( \left( \Phi(x) - 0.5 \right) \cdot \left( \Phi(y) - 0.5 \right) \cdot \phi_2(x,y|c_j^\star) \right) dx\, dy, \quad 2 \leqslant j \leqslant n. \tag{167}$$

In (167), $\phi_2(x, y | c_j^\star)$ is the probability density function of a 2-dimensional Gaussian random variable with zero mean, unit variance, and the correlation between its two elements equal to $c_j^\star$. Elements of $\mathbf{p}$ are given by

$$p_j = \begin{cases} n/(24 \cdot (\omega - 1)), & 2 \leqslant j \leqslant \omega \\ n/(24 \cdot (\omega - 1)), & n + 2 - \omega \leqslant j \leqslant n \,, \\ 0, & \text{otherwise} \end{cases} \tag{168}$$

If all elements of $\mathbf{Fc}^\star$ were nonnegative, the $\mathbf{m}'$ defined by (164) would have a rectangular power spectrum with cutoff frequency $\omega$, and $m_i' \sim \mathcal{U}(0, 1)$ for every $i \in \mathcal{I}_n$. However, in reality, such random vector does not exist. Therefore, (166) is used to replace $\mathbf{Fc}^\star$ by the closest point in $\mathbb{C}^n$ that is elementwise nonnegative. This modification expands the power spectrum of the resulting $\Phi(\mathbf{F}^{-1}\mathbf{r})$ in (164) beyond the intended one, which is corrected by mapping $\Phi(\mathbf{F}^{-1}\mathbf{r})$ onto $\mathcal{M}_\omega$ in (164).

The window "function" $\boldsymbol{w}$ in (160) and (163) is used to scale the modulator to zero smoothly at the boundaries with no effect on the remaining points:

$$w_i = \begin{cases} \sin^2\left(\frac{\pi \cdot (i-1)}{2 \cdot n_{trn}}\right), & 1 \leqslant i \leqslant n_{trn} \\ 1, & n_{trn} < i \leqslant n - n_{trn} \qquad \forall i \in \mathcal{I}. \\ \cos^2\left(\frac{\pi \cdot (i-n+n_{trn})}{2 \cdot n_{trn}}\right), & n - n_{trn} < i \leqslant n \end{cases} \tag{169}$$

In (169), $n_{trn}$ is the length of the transition window. We assumed $n_{trn} = 3 \cdot \lceil f_s/\omega \rceil$ in the present work.

For simulations discussed in Sections IV – VI of the main text, we used $f_s = 4\,\text{kHz}$. $\omega$ was set to $15\,\text{Hz}$ for *nonstationary Gaussian* modulators and $20\,\text{Hz}$ for *maximally-uniformly distributed* modulators. The cutoff frequency control parameter $\varpi$ of the AP algorithms was fixed to $30\,\text{Hz}$.

*Carriers (recovery tests)*

For signal recovery tests discussed in Section C, carriers were generated by uniformly sampling from a subset of $\mathcal{C}_d$ formed by all spike-train carriers with a fixed number of spikes $n_s$:

$$\mathcal{C}_d^{n_s} = \left\{\mathbf{x} \in \mathcal{C}_d : \left(\sum_{i=1}^n I_{\{1\}}(x_i) = n_s\right) \wedge \left(\sum_{i=1}^n I_{\{0\}}(x_i) = n - n_s\right)\right\} \tag{170}$$

The sampling was achieved by exploiting a customized version of the Metropolis-Hastings algorithm. In particular, starting with some $\mathbf{q}^{(0)} \in \mathcal{Q}$ such that

$$\mathcal{Q} = \left\{\mathbf{x} \in \mathbb{N}_0^{n*} : \left(\sum_{i=1}^{n*} x_i = n_-\right) \wedge \left(x_i < d \; \forall i \in \mathcal{I}_{n*}\right)\right\}, \tag{171}$$

where $n_- = n_s \cdot d - n$, and $n_* = \min\{n_-, n_s - 1\}$, a sequence $\mathbf{q}^{(0)}, \mathbf{q}^{(1)}, \ldots$ defined by

$$\mathbf{q}^{(i)} = \begin{cases} \mathbf{q}^{(i-1)} + z_i \cdot (\mathbf{e}^{(l_i)} - \mathbf{e}^{(k_i)}), & \text{if } \left(\mathbf{q}^{(i-1)} + z_i \cdot (\mathbf{e}^{(l_i)} - \mathbf{e}^{(k_i)})\right) \in \mathcal{Q} \\ \mathbf{q}^{(i-1)}, & \text{otherwise} \end{cases} \tag{172}$$

for $i \geqslant 1$ was generated. In (172), $\mathbf{e}^{(l)}$ is an element of $\mathbb{N}_0^{n_*}$ with its $l$-th component equal to one and the remaining components equal to zero; $l_i$ and $k_i$ are independent random numbers from a uniform distribution on $\mathcal{I}_{n_*}$; $z_i$ is a random number from a uniform distribution on the set of all integer numbers between $-\sigma$ and $+\sigma$, where $\sigma$ is a positive integer chosen to achieve the acceptance rate of the sampling scheme (172) between 0.4 and 0.6.

The iterative process generating the sequence $\mathbf{q}^{(0)}, \mathbf{q}^{(1)}, \ldots$ was terminated upon the first instance of adherence to the following equilibration criterion of the underlying Markov chain:

$$\frac{1}{n_*} \cdot \sum_{j=1}^{n_*} \left| \frac{1}{i} \cdot \sum_{l=1}^{i} q_j^{(i)} - \frac{n_-}{n_*} \right| < 0.01 \cdot \frac{n_-}{n_*}. \tag{173}$$

The corresponding $\mathbf{q}^{(i)}$ was then taken as a sample point from $\mathcal{Q}$. Next, an extended vector $\bar{\mathbf{q}} \in \mathbb{N}_0^{n_s}$ with

$$\bar{q}_j = \begin{cases} d - q_j, & \text{if } 1 \leqslant j \leqslant n_* \\ d, & n_* + 1 \leqslant j \leqslant n_s \end{cases} \tag{174}$$

was defined. Its components were then randomly permuted to produce another vector $\tilde{\mathbf{q}} \in \mathbb{N}_0^{n_s}$. The latter was used to create an $\mathbf{r} \in \mathbb{N}_+^{n_s}$ with

$$r_j = \eta + \sum_{i=1}^{j} \tilde{q}_i, \quad j \in \mathcal{I}_{n_s}, \tag{175}$$

where $\eta$ is a random integer number (the same for all $j \in \mathcal{I}_{n_s}$) from a uniform distribution on $\mathcal{I}_n$, and the summation is assumed to be modulo $\mathcal{I}_n$. Finally, the carrier $\mathbf{c}$ was generated by taking the zero element of $\mathbb{R}^n$ and setting its components whose indexes are defined by the components of $\mathbf{r}$ to one.

*Carriers (performance tests)*

For the performance, convergence, and robustness tests of the demodulation algorithm discussed in Sections IV – VI of the main text, four types of carriers were considered:

- *Nonstationary sinusoid*,

$$\mathbf{c} = \cos\big(2\pi(f\mathbf{t} + \boldsymbol{\psi})\big), \tag{176}$$

  with $f = 200\,\text{Hz}$,

$$\mathbf{t} = f_s^{-1} \cdot (0, 1, \ldots, n-1)^{\mathrm{T}}, \tag{177}$$

  and

$$\boldsymbol{\psi} = (\mathbf{g} - \min[\mathbf{g}])/\max[\mathbf{g} - \min[\mathbf{g}]], \tag{178}$$

  where $\mathbf{g} = \mathbf{P}_{\mathcal{S}_\varpi}[\mathbf{u}]$

$$u_i \sim \mathcal{U}(0, 1), \quad i \in \mathcal{I}_n, \tag{179}$$

such that $u_i \perp\!\!\!\perp u_j$ whenever $i \neq j$.

- *Nonstationary harmonic*,

$$\mathbf{c} = \sum_{l=1}^{n_f} \left( (2/3)^{l-1} \cdot \cos\left(2\pi l f(\mathbf{t} + \boldsymbol{\psi}_l) + \eta_l\right) \right), \tag{180}$$

with $f = 85\,\mathrm{Hz}$, $n_f = \lfloor f_s/(2 \cdot f) \rfloor$,

$$\eta_l \sim \mathcal{U}(0,1), \quad l \in \mathcal{I}_{n_f}, \tag{181}$$

and

$$\boldsymbol{\psi}_l = \eta_l + 0.2 \cdot (\mathbf{g} - \min[\mathbf{g}])/\max[\mathbf{g} - \min[\mathbf{g}]], \tag{182}$$

where $\mathbf{g} = \mathbf{P}_{\omega'}[\mathbf{u}]$ with $\omega' = 1\,\mathrm{Hz}$ and

$$u_i \sim \mathcal{U}(0,1), \quad i \in \mathcal{I}_n, \tag{183}$$

such that $\eta_i \perp\!\!\!\perp \eta_j$ and $u_i \perp\!\!\!\perp u_j$ whenever $i \neq j$.

- *Nonstationary spikes*,

$$\mathbf{c} = \mathbf{1} - \theta\left(\mathbf{1} - \sum_{i=1}^{n_s} \mathbf{h} * \mathbf{e}^{(r_i)}\right) \circ \left(\mathbf{1} - \sum_{i=1}^{n_s} \mathbf{h} * \mathbf{e}^{(r_i)}\right), \tag{184}$$

where $\mathbf{e}^{(r_i)}$ is the unit vector with all but the $r_i$-th of its components equal to zero, $\mathbf{h}$ is defined by

$$h_i = \begin{cases} e^{-(i-1)^2/4}, & \text{if } 1 \leqslant i \leqslant 11 \\ e^{-(n-i+1)^2/4}, & \text{if } n - 9 \leqslant i \leqslant n\,, \\ 0, & \text{otherwise} \end{cases} \tag{185}$$

and $r_1, r_2, \ldots, r_{n_s}$ is a sequence of different elements of $\mathcal{I}_n$ generated following

$$r_i = r_{i-1} + d_i + \eta, \tag{186}$$

until $r_{i-1} \leqslant n - d_i$ with $r_1 = 1$. In (186), $\eta$ is a random number from a uniform distribution on $\mathcal{I}_n$ (the same for all $i$); $d_i$ is a random number taken independently from a uniform distribution on $\mathcal{I}_{z_i}$ at each iteration, where

$$z_i = \lceil f_s/(2\varpi) \cdot (0.8 + 0.2 \cdot \sin(2\pi 5 t_i)) \rceil, \quad i \in \mathcal{I}_n. \tag{187}$$

Note that, differently from the spike-train carriers generated for the purpose of recovery tests, (184) defines sequences of finite-width spikes.

- *White-noise*,

$$c_i \sim \mathcal{U}(-1,1), \quad i \in \mathcal{I}_n, \tag{188}$$

with $c_i \perp\!\!\!\perp c_j$ whenever $i \neq j$.

**PERFORMANCE TESTS**

### J. Configurations of Demodulation Algorithms for Performance Analysis

The following configurations of the AP and LDC demodulation algorithms were used for the performance analysis in Section IV of the main text.

*Window splitting*

Segment lengths $n_{seg} = \{2^6, 2^7, \ldots, 2^{15}\}$ were examined in the case of the AP demodulation approach. In the case of the LDC demodulation method, the range of segment lengths was more limited, $n_{seg} = \{2^6, 2^7, \ldots, 2^{11}\}$, to make the simulation times feasible. In order to reduce demodulation errors stemming from the window decomposition, signals were split into segments with a particular overlap. On top of that, each segment was windowed with the Hann function [35]. For each segment length $n_{seg}$, different overlap spans were assumed: $n_{olp} = \{2^5, 2^6, \ldots, n_{seg}/2\}$. The AS-based demodulation was performed only with the original signal windows using no decomposition.

*Control parameters of the AP algorithms*

Overall, only two parameters are associated with the AP demodulation algorithms: 1) the cutoff frequency $\varpi$; and 2) the number of iterations $N_{iter}$. In all cases, we fixed $\varpi$ to 40 Hz. A set of different values of $N_{iter}$ was considered, dependent on the particular algorithm: the range from 1 to 600 for AP-B, 1 to 40 for AP-A, and 1 to 6000 for AP-P. It was made sure that the maximum iteration numbers in all these sets allow for reaching demodulation precision that is high enough not to affect the conclusions of the performance analysis.

*Control parameters of the quadratic programming solvers for the LDC approach*

In the case of the LDC approach, all elements of the weighting vector $\mathbf{w}$ [see (11)] corresponding to frequencies below the threshold $\varpi$ were assumed to be zero. The remaining elements of $\mathbf{w}$ were set to either of $\{10^2, 10^3\}$. In the case of the OSQP solver, the following control parameters were tuned:

- Linear System Solver: $\{$*Suite-Sparse-LDL, MKL-Paradiso*$\}$;
- Solution Polishing: $\{false, true\}$;
- Warm Starting: $\{false, true\}$;
- Absolute Tolerance: $\{10^{-2}, 10^{-3}, \ldots, 10^{-6}\}$;
- Relative Tolerance = Absolute Tolerance;
- Primal Infeasibility Tolerance = Absolute Tolerance;
- Dual Infeasibility Tolerance = Absolute Tolerance;
- Maximum Iteration Number: $2^{15} - 1$;

- Run Time Limit: 0.

The Gurobi solver was tried with these settings:

- Method: $\{0, 1, 2\} \equiv \{$*primal-simplex*, *dual-simplex*, *barrier*$\}$;
- Optimality Tolerance: $\{10^{-2}, 10^{-3}, \ldots, 10^{-6}\}$;
- Feasibility Tolerance = Optimality Tolerance;
- Run Time Limit: 0.

## K. IMPLEMENTATION ON A COMPUTER

All demodulation algorithms considered in the present work were implemented in C and then interfaced with MATLAB (R2018a) for a large-scale management of different instances and data analysis. The C code was compiled with GCC (v8.3) using no optimization. Evaluation of the discrete Fourier transform, used in the AS and AP approaches, relied on the FFT implementation of the Intel Math Kernel Library 2019 (update 5). The calculations were performed on a Lenovo "ThinkCentre M910t Tower" desktop computer with an Intel Core i7-7700 processor and Linux Ubuntu 16.04 operating system. In order to minimize the influence of other operating system processes on the benchmarking results, one of the four CPU cores was dedicated exclusively to the execution of the demodulation program. This was achieved by using the "isolplus" option of the kernel scheduler. All computations were done in single-thread mode.

## ADDITIONAL RESULTS

### L. ERRORS OF CARRIER ESTIMATES

Here, we discuss findings from the performance and robustness analyses of demodulation algorithms (introduced in Sections IV and VI of the main text) in terms of carrier estimates. Analogous to the modulator recovery error $E_m$, we use $E_c = \|\mathbf{c} - \hat{\mathbf{c}}\|_2 / \|\mathbf{c}\|_2$ next.
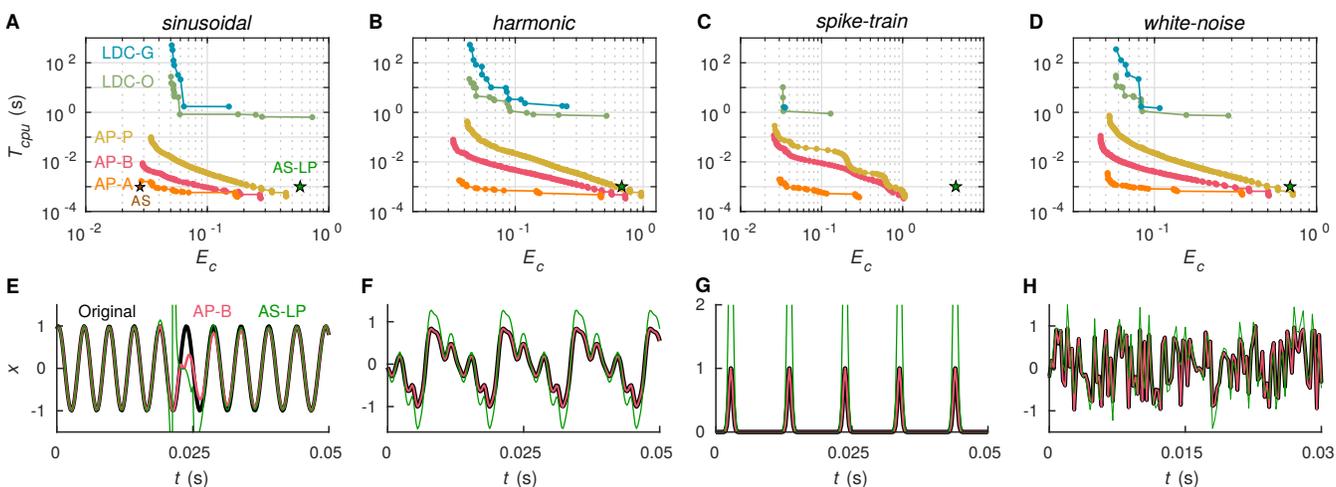
*Performance tests*



Fig. 13.  Performance evaluation. **A–D**: Pareto fronts in the $(E_c, T_{\text{cpu}})$ plane for different demodulation algorithms applied to the four different types of test signals when window splitting is used. The green and brown stars mark the results of, respectively, the AS-LP and AS methods. **E–H**: Examples of the original carriers considered in the present work (black) and their estimates obtained by using the AP-B (red) and AS-LP (green) algorithms.

The high precision of modulator estimates (see Section IV-C in the main text) and the boundedness of $\hat{c}_i$ between $-1$ and $1$ (see Section VII in the main text) predetermine good quality carrier estimates provided by the AP approach.[20] This view is evidenced by the Pareto fronts in the $(E_c, T_{\text{cpu}})$ plane shown in Fig. 13 A–D and comparisons of exemplary $\mathbf{c}$ and $\hat{\mathbf{c}}$ in Fig. 13 E–H. We observed noticeable discrepancies between $\mathbf{c}$ and $\hat{\mathbf{c}}$ only locally, around points with modulator levels very close to zero (see the signal segment at $t = 0.025$ in Fig. 13 E). That finding is explained by the fact that, in our case, $\hat{c}_i - c_i = s_i \cdot (\hat{m}_i^{-1} - m_i^{-1})$. Hence, even small differences between $m_i$ and $\hat{m}_i$ on the absolute scale can give notable differences between $c_i$ and $\hat{c}_i$ if $m_i/\hat{m}_i \gg 1$. However, the condition $\hat{m}_i \geqslant s_i$ ( $\implies |c_i| \leqslant 1$ ) inherent to all $\hat{\mathbf{c}}$ obtained by the AP algorithms assures tolerable distortions of the carrier estimates even around the points with vanishingly small $m_i$.

---

[20]Remember that $c_i/\hat{c}_i = m_i/\hat{m}_i$ in our case.

The situation is different in the case of the AS-LP approach. Then, $|s_i/\hat{m}_i| \gg 1$, i.e., $|\hat{c}_i| \gg 1$, are possible (see the signal segment at $t = 0.025$ in Fig. 13 E). These divergences noticeably increase the recovery errors $E_c$ even for sinusoidal signals that AS-LP recovers sufficiently well outside the segments of vanishing $m_i$ (see Fig. 13 A, E). Globally-inaccurate recovery of other carrier types demonstrated by the AS-LP (see Fig. 13 B–D, F–H) is predetermined by inaccurate estimates of the modulators (see Section IV-C in the main text). We note that, for sinusoidal carriers, appropriate estimates $\hat{c}$ can be obtained by using the original AS method instead of the AS-LP (Fig. 13 A). However, the AS gives very erroneous modulator estimates $\hat{m}$, and hence, does not improve carrier predictions $\hat{c}$ considerably, for other types of signals (data not shown).
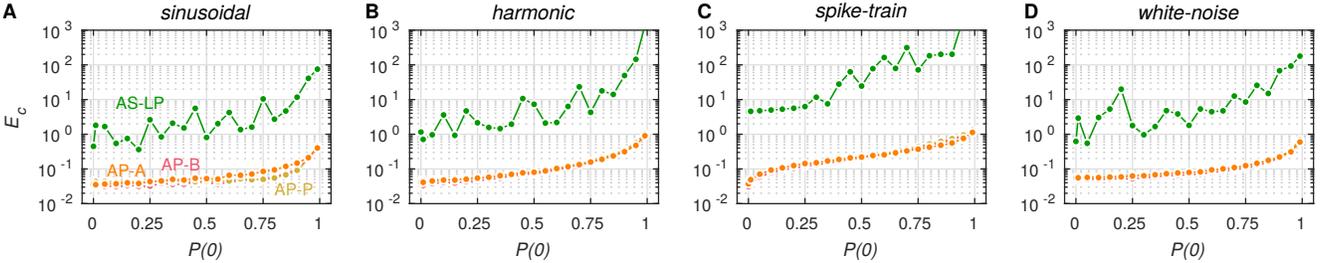
*Robustness tests*



Fig. 14. Robustness evaluation. **A–D**: Dependence of the carrier recovery error $E_c$ on $P(0)$ (the probability of missing points) for the four types of test signals and different AP algorithms at $\epsilon_{tol} = 10^{-4}$ (color coding).

Fig. 14 shows demodulation results of the AS-LP, AP-B, AP-A, and AP-P algorithms in terms of carrier recovery error for test signals corrupted by the multiplicative Bernoulli-$\{0, 1\}$ noise (see Section VI). The obtained $E_c$ vs. $P(0)$ relations for the AP algorithms are analogous to their counterparts $E_m$ vs. $P(0)$ shown in Fig. 5. In contrast, the AS-LP is even more inferior to the AP approach in terms of carrier reconstruction (Fig. 14) than it is in terms of modulator recovery (Fig. 5 A–D). This fact is explained by the presence of the $|\hat{c}_i| \gg 1$ divergences discussed above and illustrated by Fig. 13 E.

## M. CONVERGENCE RATES AT DIFFERENT SIGNAL LENGTHS

The convergence results shown in Fig. 4 of the main text represent only signals with the length $n$ fixed to $2^{15}$ sample points. Therefore, we performed additional simulations with different $n$ values to test the impact of this parameter on the rate of the iterative process. We found no clear dependence of $N_{iter}$ required to achieve a particular infeasibility error value $\epsilon$ on $n$, except transitional changes due to diminishing contributions of the boundary effects in some cases (see Fig. 15 $C_2 - C_2$). These findings suggest that the increased $T_{cpu}$ of demodulation for longer $n$ is primarily determined by the increased computational demands of single projections (see Section III-D in the main text).
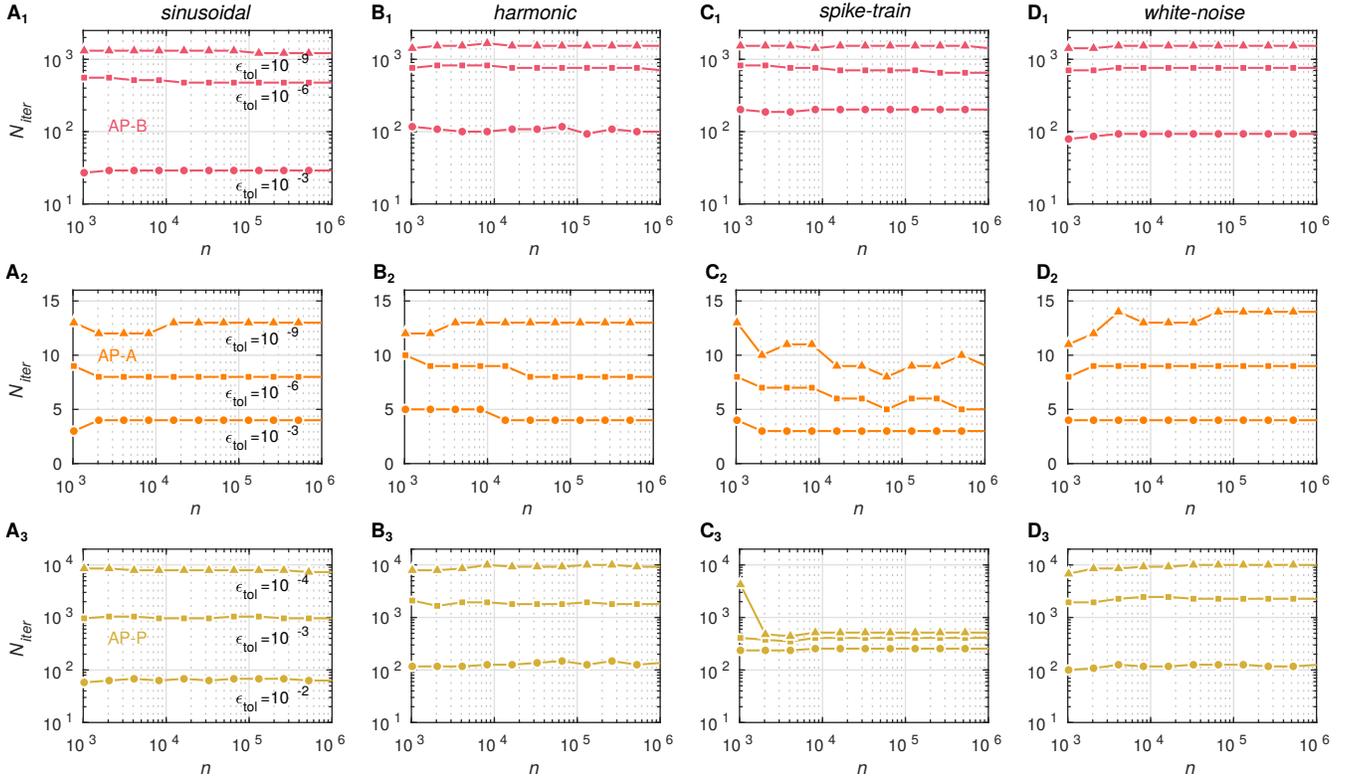
Fig. 15. Convergence analysis of the AP algorithms at different $n$ and fixed $f_s = 4\,\text{kHz}$. $\mathbf{A_1}$–$\mathbf{D_1}$: Dependence of the iteration number $N_{iter}$ necessary to reach a specific infeasibility error $\epsilon$ on the length of the signal for the AP-B algorithm applied to four different classes of test signals. Filled circles, rectangles, and triangles label curves for $\epsilon$ levels of, respectively, $10^{-3}$, $10^{-6}$, and $10^{-9}$. $\mathbf{A_2}$–$\mathbf{D_2}$: The same as $\mathrm{A_1}$–$\mathrm{D_1}$ but shown for the AP-A algorithm. $\mathbf{A_3}$–$\mathbf{D_3}$: The same as $\mathrm{A_1}$–$\mathrm{D_1}$ but shown for the AP-P algorithm and different levels of the infeasibility error $\epsilon$.
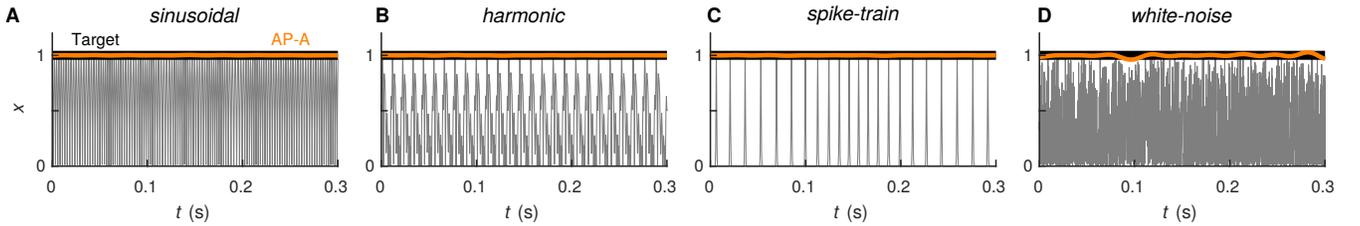


Fig. 16. Repeated demodulation of inferred carriers. Gray – sinusoidal (**A**), harmonic (**B**), spike-train (**C**), and white-noise (**D**) carriers inferred by demodulating test signals shown in Fig. 2 A–D with the AP-A algorithm. Black – target repeated modulators following from the assumption that the inferred carriers are fully demodulated. Orange – the real repeated modulators obtained with AP-A.

## N. REPEATED DEMODULATION

Fig. 16 shows typical results of redemodulation of carriers inferred from the four types of synthetic test signals considered in the present work by using the AP-A algorithm. Redemodulation of the carriers returns the identity

modulator to a very good approximation ($E_m \leqslant 10^{-2}$), implying a nearly complete separation of the modulator and carrier information in the first step, as discussed in Section VII of the main text.

## O. DEMODULATION OF SPEECH SIGNALS

In this section, we present results of additional simulations used to support the statements about the suitability of the AP approach to demodulate wideband speech signals in Section VIII of the main text.

*Synthetic modulators*

The first question that we considered was up to which values of the cutoff frequency $\omega$ the natural speech carriers meet the recovery condition $\lceil n/d \rceil \geqslant 2\omega - 1$. As mentioned in the main text, these carriers are of quasi-random and quasi-harmonic origins, possibly featuring frequency glides. It is well known that typical fundamental frequencies ($f_0$) of male and female speaker voice are, respectively, 120 and 210 Hz (see [36] and Table 1 therein). Hence, at least for the harmonic components, the condition $\lceil n/d \rceil \geqslant 2\omega - 1$ is expected to be satisfied with $\omega \leqslant 60$ Hz. That is more than sufficient for appropriate demodulation, assuming that strictly all spectral energy of the speech amplitude modulator is located below 20 Hz (see Section VIII-A). The situation with the quasi-random and mixed components of speech signals is less certain and must be tested numerically.

To this end, we considered four speech-carriers generated by a male speaker uttering prolonged ($\sim$ 1 s) [ɑ:], [ʃ], and [z], as well as vocable [waʊ] (see Fig. 17 $A_1 - A_4$) at $f_0 = 120$ Hz without noticeable amplitude modulation. The [ɑ:] is predominantly harmonic, [ʃ] is quasi-random, [z] has a mixed wave-shape, and [waʊ] is quasi-harmonic but with upward and downward frequency glides. These characteristics are further revealed by the power-spectral-density (PSD) plots for the [ɑ:], [ʃ], and [z] (see Fig. 17 $B_1 - B_3$), and a spectrogram for the [waʊ] (Fig. 17 $B_4$). We then created test signals as products of mentioned carriers and maximally-uniformly distributed synthetic modulators (see Section I) with $\omega = 25$ Hz and $f_s = 44.1$ kHz.

Demodulation of the considered test signals with the AP-A algorithm allowed us to achieve high-accuracy modulator recovery, as shown in Fig. 17 $C_1 - C_4$ (note the insets with $E_m$ and $E_c$ values there). In more detail, we found that, for the [ɑ:], [ʃ], [z], and [waʊ], the distances between carrier points with absolute values of, respectively, $\geqslant 0.99$, $\geqslant 0.95$, $\geqslant 0.93$, and $\geqslant 0.93$ were below that needed to satisfy $\lceil n/d \rceil \geqslant 2\omega - 1$ at $\omega = 20$ Hz. For the [ʃ], [z], and [waʊ], 80 % of the required points were $\geqslant 0.99$. The demodulation quality remained reasonably good when $\omega$ was increased to 50 Hz, resulting in $E_m$ values of $1 \cdot 10^{-2}$, $5 \cdot 10^{-2}$, $5 \cdot 10^{-2}$, and $5 \cdot 10^{-2}$ for, respectively, [ɑ:], [ʃ], [z], and [waʊ] carriers.

*Natural modulators*

Next, we aimed to clarify whether the AP approach can properly separate **m** and **c** of speech signals using the dynamic range compression discussed in Section VIII-B of the main text. For this purpose, we considered
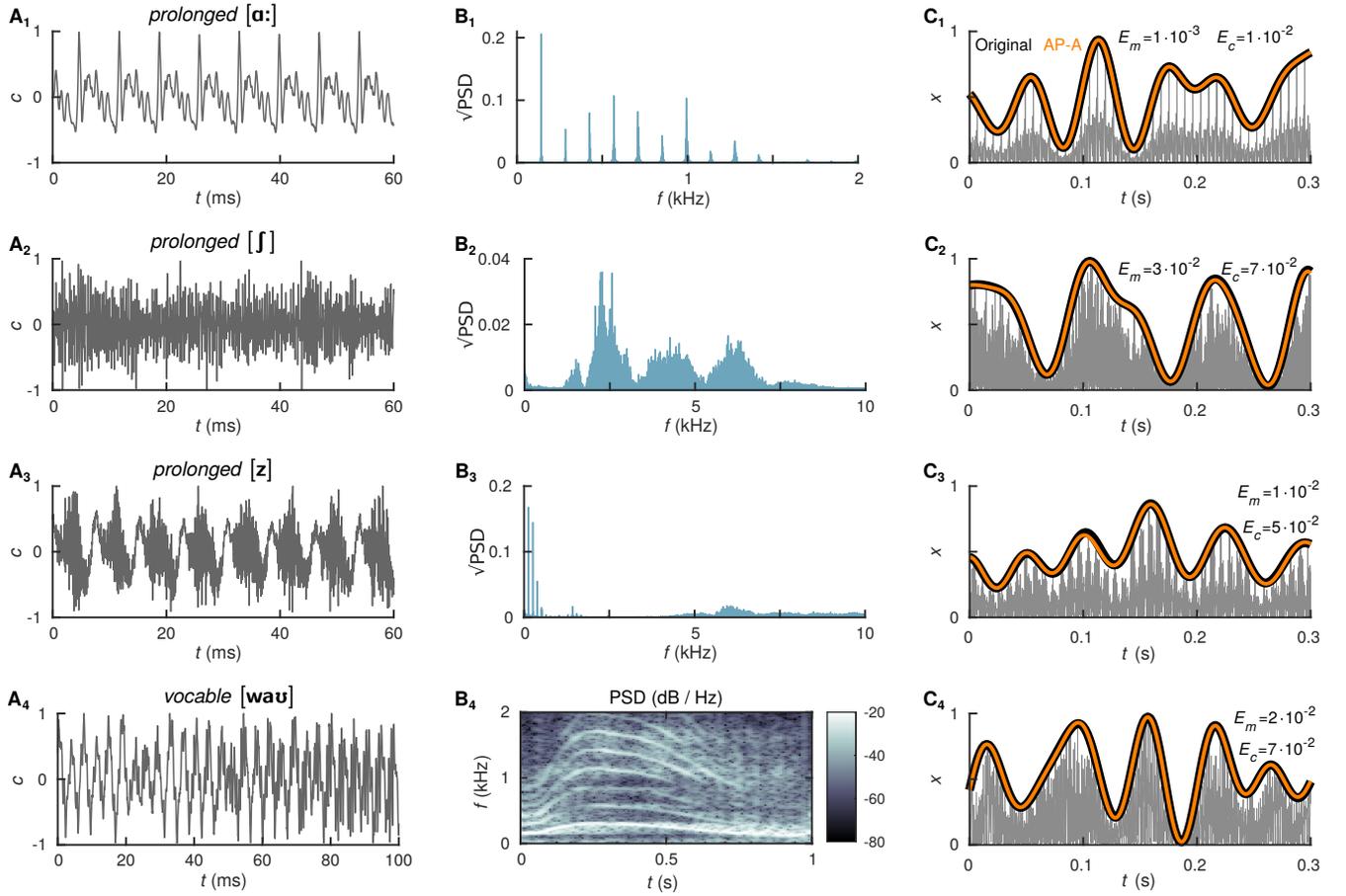
Fig. 17. Demodulation of signals built of synthetic modulators and natural speech carriers. $A_1$–$A_4$: fragments of carrier-signals of prolonged steady [ɑ:] ($A_1$), [ʃ] ($A_2$), and [z] ($A_3$), as well as an extended vocable [waʊ] ($A_4$). $B_1$–$B_3$: periodogram estimators of the power spectral densities of the carriers illustrated in $A_1$–$A_3$; $B_4$: spectrogram of the vocable [waʊ]. $C_1$–$C_4$: exemplary segments of amplitude-modulated carriers from panels $A_1$–$A_4$, their modulators (black), and modulator estimates obtained with the AP-A algorithm. Insets of panels $C_1$–$C_4$ display values of the modulator and carrier recovery errors.

synthetic time series built of the four natural carriers [ɑ:], [ʃ], [z], and [waʊ] introduced above and modulator estimate $\hat{\mathbf{m}}^*$ of an utterance *"...protein which forms p ..."* from Section VIII-B of the main text. We then demodulated the resulting signals by using the AP-A algorithm combined with the dynamic range compression. The obtained estimates $\hat{\hat{\mathbf{m}}}^*$ were in good agreement with the original, as shown in Fig. 18 (note the insets with $E_m$ and $E_c$ values there). The AP-B, AP-P, and LDC algorithms returned very similar modulator estimates but needed much longer computing times that mirror the performance analysis presented in Section IV-C of the main text.
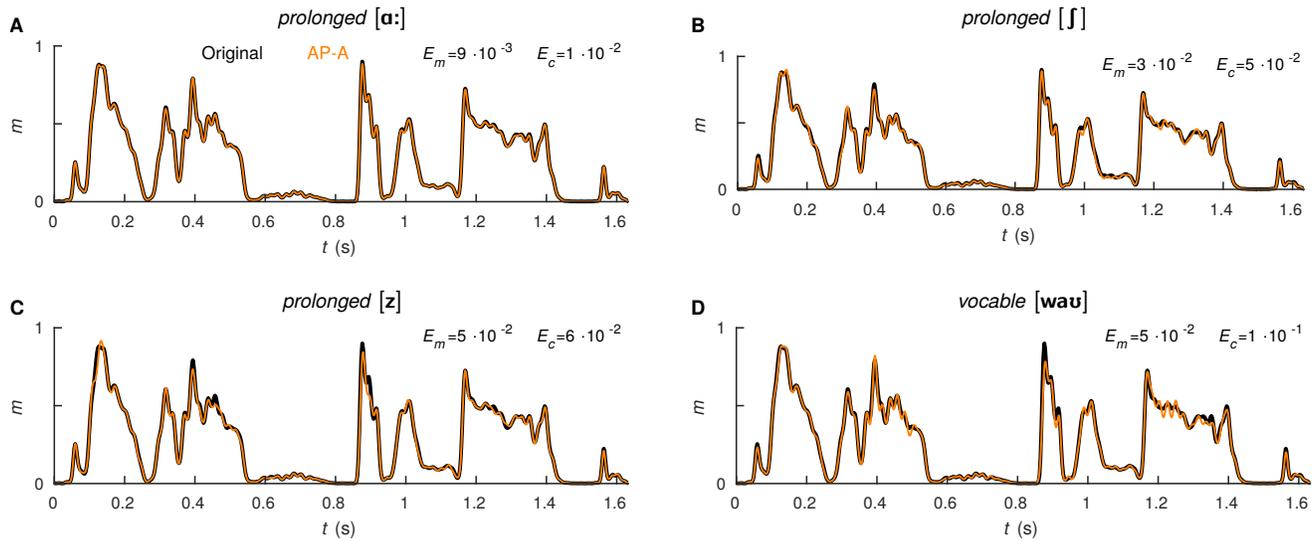
Fig. 18. Demodulation of signals built of natural speech modulators and carriers. **A–D** compares original modulators $\hat{\mathbf{m}}^*$ (black) and their estimates $\hat{\hat{\mathbf{m}}}^*$ obtained by using the AP-A algorithm combined with the dynamic range compression (orange) for, respectively, [ɑ:], [ʃ], [z], and [waʊ] carriers. Insets display values of the modulator and carrier recovery errors.

# REFERENCES

[1] F. Zhang, *Matrix Theory*. Springer, 2011.

[2] A. N. Kolmogorov and S. V. Fomin, *Introductory real analysis*. Courier Corporation, 1975.

[3] E. Çinlar and R. J. Vanderbei, *Real and Convex Analysis*. Springer, 2013.

[4] L. M. Bregman, "The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming," *USSR Comput. Math. & Math. Phys.*, vol. 7, pp. 200–217, 1967.

[5] L. G. Gubin, B. T. Polyak, and E. V. Raik, "The method of projections for finding the common point of convex sets," *USSR Comput. Math. & Math. Phys.*, vol. 7, pp. 1–24, 1967.

[6] J. P. Boyle and R. L. Dykstra, "A method for finding projections onto the intersection of convex sets in Hilbert spaces," in *Advances in Order Restricted Statistical Inference*, ser. Lecture Notes in Statistics. Springer, 1986, pp. 28–47.

[7] V. F. Humphrey, "Nonlinear propagation in ultrasonic fields: measurements, modelling and harmonic imaging," *Ultrasonics*, vol. 38, pp. 267–272, 2000.

[8] F. A. Duck, "Nonlinear acoustics in diagnostic ultrasound," *Ultrasound Med. Biol.*, vol. 28, pp. 1–18, 2002.

[9] L. Sornmo and P. Laguna, *Bioelectrical Signal Processing in Cardiac and Neurological Applications*. Academic Press, 2005.

[10] G. L. Gottlieb and G. C. Agarwal, "Filtering of electromyographic signals," *Am. J. Phys. Med. Rehab.*, vol. 49, p. 142, 1970.

[11] T. J. Roberts and A. M. Gabaldón, "Interpreting muscle function from EMG: lessons learned from direct measurements of muscle force," *Integr. Comp. Biol.*, vol. 48, pp. 312–320, 2008.

[12] "UCI Machine Learning Repository: sEMG for Basic Hand movements Data Set." [Online]. Available: https://archive.ics.uci.edu/ml/datasets/sEMG+for+Basic+Hand+movements

[13] C. Sapsanis, "Recognition of basic hand movements using Electromyography," Diploma Thesis, University of Patras, 2013, arXiv: 1810.10062.

[14] J. E. Shoup and L. L. Pfeifer, "Acoustic characteristics of speech sounds," in *Contemporary Issues in Experimental Phonetics*. Academic Press, 1976, pp. 171–224.

[15] R. E. Turner, "Statistical models for natural sounds," Ph.D. dissertation, University College London, 2010. [Online]. Available: http://discovery.ucl.ac.uk/19231/

[16] A. Keitel, J. Gross, and C. Kayser, "Perceptually relevant speech tracking in auditory and motor cortex reflects distinct linguistic features," *PLOS Biol.*, vol. 16, p. e2004473, 2018.

[17] B. E. D. Kingsbury, N. Morgan, and S. Greenberg, "Robust speech recognition using the modulation spectrogram," *Speech Commun.*, vol. 25, pp. 117–132, 1998.

[18] S. Wu, T. H. Falk, and W.-Y. Chan, "Automatic speech emotion recognition using modulation spectral features," *Speech Commun.*, vol. 53, pp. 768–785, 2011.

[19] B. Lee and K.-H. Cho, "Brain-inspired speech segmentation for automatic speech recognition using the speech envelope as a temporal reference," *Sci. Rep.*, vol. 6, p. 37647, 2016.

[20] B. S. Wilson, C. C. Finley, D. T. Lawson, R. D. Wolford, D. K. Eddington, and W. M. Rabinowitz, "Better speech recognition with cochlear implants," *Nature*, vol. 352, pp. 236–238, 1991.

[21] F.-G. Zeng, K. Nie, G. S. Stickney, Y.-Y. Kong, M. Vongphoe, A. Bhargave, C. Wei, and K. Cao, "Speech recognition with amplitude and frequency modulations," *PNAS*, vol. 102, pp. 2293–2298, 2005.

[22] R. V. Shannon, F.-G. Zeng, V. Kamath, J. Wygonski, and M. Ekelid, "Speech recognition with primarily temporal cues," *Science*, vol. 270, pp. 303–304, 1995.

[23] Z. M. Smith, B. Delgutte, and A. J. Oxenham, "Chimaeric sounds reveal dichotomies in auditory perception," *Nature*, vol. 416, pp. 87–90, 2002.

[24] P. X. Joris, C. E. Schreiner, and A. Rees, "Neural processing of amplitude-modulated sounds," *Physiol. Rev.*, vol. 84, pp. 541–577, 2004.

[25] U. Goswami, "Speech rhythm and language acquisition: An amplitude modulation phase hierarchy perspective," *Ann. N. Y. Acad. Sci.*, vol. 1453, pp. 67–78, 2019.

[26] A. A. Sarkady, R. R. Clark, and R. Williams, "Computer analysis techniques for phonocardiogram diagnosis," *Comput. Biomed. Res.*, vol. 9, pp. 349–363, 1976.

[27] D. Gill, N. Gavrieli, and N. Intrator, "Detection and identification of heart sounds using homomorphic envelogram and self-organizing probabilistic model," in *Computers in Cardiology, 2005*, 2005, pp. 957–960.

[28] B. Oliver, J. Pierce, and C. Shannon, "The philosophy of PCM," *Proc. IRE*, vol. 36, pp. 1324–1331, 1948.

[29] F. Rieke, D. Warland, R. de Ruyter van Steveninck, and W. Bialek, *Spikes: exploring the neural code*.   A Bradford Book, 1999.

[30] J. Felblinger and C. Boesch, "Amplitude demodulation of the electrocardiogram signal (ECG) for respiration monitoring and compensation during MR examinations," *Magn. Reson. Med.*, vol. 38, pp. 129–136, 1997.

[31] G. D. Gargiulo, P. Bifulco, M. Cesarelli, A. L. McEwan, H. Moeinzadeh, A. O'Loughlin, I. M. Shugman, J. C. Tapson, and A. Thiagalingam, "On the "Zero of potential of the electric field produced by the heart beat". A machine capable of estimating this underlying persistent error in electrocardiography," *Machines*, vol. 4, p. 18, 2016.

[32] A. L. Goldberger, L. A. N. Amaral, L. Glass, J. M. Hausdorff, P. C. Ivanov, R. G. Mark, J. E. Mietus, G. B. Moody, C.-K. Peng, and H. E. S. H. Eugene, "PhysioBank, PhysioToolkit, and PhysioNet," *Circulation*, vol. 101, pp. e215–e220, 2000.

[33] L. Wasserman, *All of statistics: A concise course in statistical inference*.   Springer, 2004.

[34] M. C. Cario and B. L. Nelson, "Modeling and generating random vectors with arbitrary marginal distributions and correlation matrix," Department of Industrial Engineering and Management Sciences, Northwestern University, Tech. Rep., 1997.

[35] G. Sell and M. Slaney, "Solving demodulation as an optimization problem," *IEEE Audio, Speech, Language Process.*, vol. 18, pp. 2051–2066, 2010.

[36] H. Traunmüller and A. Eriksson, "The frequency range of the voice fundamental in the speech of male and female adults," Department of Linguistics, University of Stockholm, Tech. Rep., 1994. [Online]. Available: http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.64.5133