

This is the accepted manuscript made available via CHORUS. The article has been published as:

Machine Learning for Predictive Estimation of Qubit Dynamics Subject to Dephasing

Riddhi Swaroop Gupta and Michael J. Biercuk

Phys. Rev. Applied **9**, 064042 — Published 27 June 2018

DOI: [10.1103/PhysRevApplied.9.064042](https://doi.org/10.1103/PhysRevApplied.9.064042)

Machine learning for predictive estimation of qubit dynamics subject to dephasing

Riddhi Swaroop Gupta* and Michael J. Biercuk

*ARC Centre of Excellence for Engineered Quantum Systems, School of Physics,
The University of Sydney, New South Wales 2006, Australia*

Decoherence remains a major challenge in quantum computing hardware and a variety of physical-layer controls provide opportunities to mitigate the impact of this phenomenon through feedback and feedforward control. In this work, we compare a variety of machine learning algorithms derived from diverse fields for the task of state estimation (retrodiction) and forward prediction of future qubit state evolution for a single qubit subject to classical, non-Markovian dephasing. Our approaches involve the construction of a dynamical model capturing qubit dynamics via autoregressive or Fourier-type protocols using only a historical record of projective measurements. A detailed comparison of achievable prediction horizons, model robustness, and measurement-noise-filtering capabilities for Kalman Filters (KF) and Gaussian Process Regression (GPR) algorithms is provided. We demonstrate superior performance from the autoregressive KF relative to Fourier-based KF approaches and focus on the role of filter optimization in achieving suitable performance. Finally, we examine several realizations of GPR using different kernels and discover that these approaches are generally not suitable for forward prediction. We highlight the linkages between predictive performance and kernel structure, and identify ways in which forward predictions are susceptible to numerical artefacts.

I. INTRODUCTION

In predictive estimation, a dynamically evolving system is observed and any temporal correlations encoded in the observations are used to predict the future state of the system. This generic problem is well studied in diverse fields such as engineering, econometrics, meteorology, and seismology [1–5], and is addressed in the control-theoretic literature as a form of filtering. Applying these approaches to state estimation on qubits is complicated by a variety of factors; dominant among these is the violation of the assumption of linearity inherent in most filtering applications as qubit states are formally bilinear. The case of an idling, or freely evolving qubit subject to dephasing is more complicated still, as an a priori model of system evolution suitable for implementation within standard filtering algorithms will not in general be available.

Fortunately there are many lessons to learn from classical control, even in the presence of such complications. For classical systems, machine learning techniques have enabled state tracking, control, and forecasting for highly non-linear and noisy dynamical trajectories or complex measurement protocols (e.g. [6–10]). These demonstrations move far beyond the simplified assumptions underlying many basic filtering tasks such as linear dynamics and white (uncorrelated) noise processes. For instance, so-called particle-based Bayesian frameworks (e.g. particle filtering, unscented or sigma-point filtering) allow state estimation and tracking in the presence of non-linearities in system dynamics or measurement protocols [11]. Further extensions approach the needs of a stochastically evolving system; recently, an ensemble of so-called unscented Kalman filters, named after the

underlying mathematical transformation, demonstrated state estimation and forward predictions for chaotic, non-linear systems in the absence of a prescribed model [10]. For non-chaotic, multi-component stationary random signals, other algorithmic approaches have been particularly useful for tracking instantaneous frequency and phase information, [12, 13], enabling short-run forecasting.

In the field of quantum control, work has begun to incorporate the additional challenges faced when considering state estimation on qubits, notably quantum-state collapse under projective measurement. Under such circumstances, in which the measurement backaction strongly influences the quantum state (in contrast with the classical case), it is not straightforward to extend machine learning predictive estimation techniques. Work to date has approached the analysis of projective measurement records on qubits as pattern recognition or image reconstruction problems, for example, in characterising the initial or final state of quantum system (e.g. [14–16]) or reconstructing the historical evolution of a quantum system based on large measurement records (e.g. [17–22]). In adaptive or sequential Bayesian learning applications, a projective measurement protocol may be designed or adaptively manipulated to efficiently yield noise-filtered information about a quantum system (e.g. [23–26]).

The demonstrations above typically assume the object of interest is either static, or stochastically evolves in a manner which is dynamically uncorrelated in time (white) as measurement protocols are repeated. This simplifying assumption falls well short of typical laboratory based experiments where noise processes are frequently correlated in time, and evolution may also occur rapidly relative to a measurement protocol. In such a circumstance, further complexity is introduced as the Markov condition commonly assumed in Bayesian learning frameworks [11] is immediately violated. Even in the

* rgup9526@uni.sydney.edu.au

classical case, the problem of designing an appropriate representation of non-Markovian dynamics in Bayesian learning frameworks is an active area of research (e.g [27]). Hence, the canonical real-time tracking and prediction problem - where a non-linear, stochastic trajectory of a system is tracked using noisy measurements and short-run forecasts are made - is under-explored for quantum systems with projective measurements.

In this manuscript, we develop and explore a broad class of predictive estimation algorithms allowing us to track a qubit state undergoing *stochastic but temporally correlated* evolution using a record of projective measurements, and forecast its future evolution. Our approaches employ machine learning algorithms to extract temporal correlations from the measurement record and use this information to build an effective dynamical model of the system's evolution. We design a deterministic protocol to correlate Markovian processes such that a certain general class of non-Markovian dynamics can be approximately tracked without violating the assumptions of a machine learning protocol, based on the theoretically accessible and computationally efficient frameworks of Kalman Filtering (KF) and Gaussian Process Regression (GPR). Both frameworks provide a mechanism by which temporal correlations (equally, dynamics) are encoded into an algorithm's structure such that projection of data-sets onto this structure enables meaningful learning, white-noise filtering, and effective forward prediction. We perform numerical simulations to test the effectiveness of these algorithms in maximizing the prediction horizon under various conditions, and quantify the role of the measurement sampling rate relative to the noise dynamics in defining the prediction horizon. Simulations incorporate a variety of measurement models, including pre-processed data yielding a continuous measurement outcome and discretised outcomes commonly associated with single-shot projective qubit measurements. We find that in most circumstances an autoregressive Kalman framework yields the best performance, providing model-robust forward prediction horizons and effective filtering of measurement noise. Finally, we demonstrate that standard GPR-based protocols employing a variety of kernels, while effective for the problem of filtering (fitting) a measurement record, are not suitable for real-time forecasting beyond the measurement record.

In what follows, we describe in detail the physical setting for our problem in Section II and explain how this leads to a specific choice of algorithm which may be deployed for the task of tracking non-Markovian state dynamics in the absence of a dynamical model for system evolution. We provide an overview of the central GPR and KF frameworks in Section III, and we specify a series of algorithms under consideration in this paper tailored to different measurement processes. For pre-processed measurement records, we consider four algorithmic approaches: a Least Squares Filter (LSF) from [28]; an Autoregressive Kalman Filter (AKF); a so-called Liska Kalman Filter from [29] adapted for a Fixed oscillator

Basis (LKFFB); and a suitably designed GPR learning protocol. For binary measurement outcomes, we extend the AKF to a Quantised Kalman Filter (QKF). In Section IV A, we present optimisation procedures for tuning all algorithms. Numerical investigations of algorithmic performance are presented in Section IV and a comparative analysis of all algorithms is provided in Section V.

II. PHYSICAL SETTING

Our physical setting considers a sequence of projective measurements performed on a qubit. Each projective measurement yields a 0 or 1 outcome representing the state of the qubit. The qubit is then reset, and the exact procedure is repeated. By considering a qubit state initialized in a superposition of the measurement basis (for us, Pauli $\hat{\sigma}_z$ eigenstates), we gain access to a direct probe of qubit phase evolution. If, for instance, no dephasing is present, then the probability of obtaining a binary outcome remains static in time as sequential qubit measurements are performed. If slowly drifting environmental dephasing is present, then the probability of obtaining a given binary outcome also drifts stochastically. In essence, the qubit probes dephasing noise and our procedure encodes a continuous-time non-Markovian dephasing process into time-stamped, discrete binary samples through the nonlinear projective measurement, carrying the underlying correlations in the noise. It is this series of measurements which we seek to process in our algorithmic approaches to qubit state tracking and prediction.

Formally, an arbitrary environmental dephasing process manifests as time-dependent stochastic detuning, $\delta\omega(t)$, between the qubit frequency and the system master clock. This detuning is an experimentally measurable quantity in a Ramsey protocol, as shown schematically in Fig. 1 (a). A non-zero detuning over measurement period τ (starting from $t = 0$) induces a stochastic relative phase accumulation (in the rotating frame) for a qubit superposition state as $|0\rangle + e^{-if(0,\tau)}|1\rangle$ between qubit basis states. The accumulated $f(0,\tau)$ at the end of a single Ramsey experiment is mapped to a probability of obtaining a particular outcome in the measurement basis via the form of the Ramsey sequence.

In a sequence of n Ramsey measurements spaced Δt apart with a fixed duration, τ , the change in the statistics of measured outcomes over this measurement record depends solely on the dephasing $\delta\omega(t)$. We assume that the measurement action over τ is much faster than the temporal dynamics of the dephasing process, and $\Delta t \gtrsim \tau$. The resulting measurement record is a set of binary outcomes, $\{d_n\}$, determined probabilistically from n true stochastic qubit phases, $f := \{f_n\}$. Here the accumulated phase in each Ramsey experiment, $f(n\Delta t, n\Delta t + \tau) \equiv \int_{n\Delta t}^{n\Delta t + \tau} \delta\omega(t') dt'$ and we use the shorthand $f(n\Delta t, n\Delta t + \tau) \equiv f_n$. We define the statistical likelihood for observing a single shot, d_n , using Born's rule [30]:

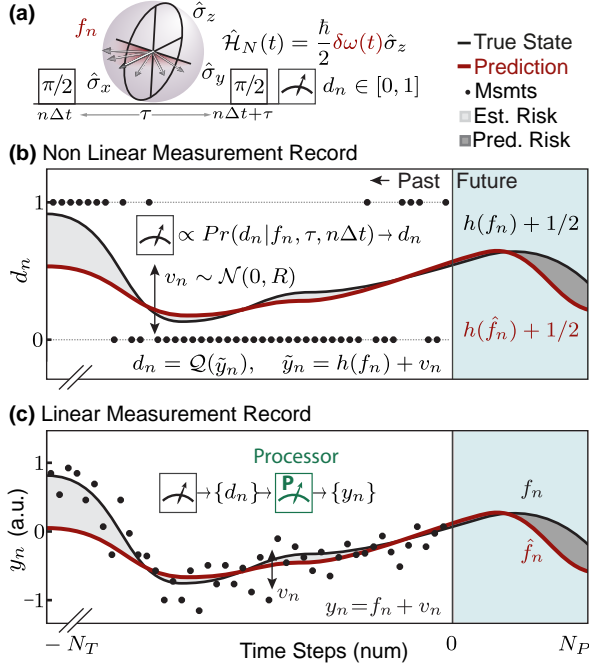


FIG. 1. (a) A Ramsey experiment at $t = n\Delta t$ with fixed wait time τ and time-steps, n , spaced $\Delta t > \tau$ apart. A $\pi/2$ pulse rotates qubit state to super-position of $|d\rangle$ states, $d \in \{0, 1\}$; qubit evolves via $\hat{H}_N(t)$ accumulating relative stochastic f_n , for non-zero environmental dephasing $\delta\omega(t)$. Jittering arrows depict potential qubit state vectors permitted for (unknown) random f_n . Qubit state is measured as $d_n = d$ in $\hat{\sigma}_z$ basis after a second $\pi/2$ rotation. (b) Black dots depict $\{d_n\}$ against time steps, n ; data collection stops at $n = 0$ separating past state estimation from future prediction [blue region]. Black solid line shows true qubit state likelihood $\propto h(f_n)$; and red solid line shows state estimate (prediction) for $n < 0$ ($n > 0$). A prediction horizon is $n < n^* \in [0, N_P]$ for which dark-grey region between red and black lines is minimised (Bayes prediction risk) relative to predicting the mean of dephasing noise; algorithmic tuning occurs by minimising light-grey region (Bayes state estimation risk). \mathcal{Q} quantises black line into noisy qubit measurements, d_n , under Gaussian uncertainty v_n . (c) Single shot outcomes in (b) are pre-processed to yield noisy measurements $\{y_n\}$ [black dots]; y_n is linear in f_n and v_n represents additive white Gaussian measurement noise.

$$Pr(d_n = d|f_n, \tau, n\Delta t) = \begin{cases} \cos^2(\frac{f_n}{2}) & \text{for } d = 1 \\ \sin^2(\frac{f_n}{2}) & \text{for } d = 0 \end{cases} \quad (1)$$

The notation $Pr(d_n|f_n, \tau, n\Delta t)$ refers to the conditional probability of obtaining measurement outcome d_n given a true stochastic phase, f_n , accumulated over τ , beginning at time $t = n\Delta t$. In the noiseless case, $Pr(d_n = 1|f_n, \tau, n\Delta t) = 1$, $\forall n$, such that a qubit exhibits no additional phase accumulation due to environmental dephasing. Following a single measurement the qubit state is reset, but the dephasing noise correlations manifest

again via Born's rule for another random value of the bias at time-step $n + 1$. A detailed discussion of Eq. (1) can be found in Appendix A.

The action of measurement, expressed as $h(f_n)$, is given by $Pr(d_n = d|f_n, \tau, n\Delta t) \equiv \frac{1}{2} - (-1)^d h(f_n)$ and is depicted in Fig. 1(b) as a probability of seeing the qubit in the $d = 1$ state. We begin by describing here a 'raw' non-linear measurement record, $\{d_n\}$ where each d_n [black dots] corresponds to a binary outcome derived from a single projective measurement on a qubit. The sequence $\{d_n\}$ can be treated as a sequence of biased coin flips, where the underlying bias of the coin is a non-Markovian, discrete-time process and the value of the bias is given by Eq. (1) at each n . The non-linearity of the measurement, $h(f_n)$, is defined with respect to f_n where Eq. (1) is interpreted as a non-linear measurement action for Bayesian learning frameworks.

This data series is contrasted with a linear measurement record, $\{y_n\}$, depicted in Fig. 1(c). Each value y_n is derived from the sum of a true qubit phase, f_n , and Gaussian white measurement noise, v_n . The sequence $\{y_n\}$ is generated by pre-processing raw binary measurements, $\{d_n\}$ via a range of experimental techniques subject to a separation of timescales such that $\sim \tau$ is much faster than drift of $\delta\omega(t)$. In the most common case, one performs M runs of the experiment over which $\delta\omega(t)$ is approximately constant, giving an estimate of f_n at $t = n\Delta t$ using averaging, a Bayesian scheme, or Fourier analysis. A more complex linearization protocol involves the use of low-pass or decimation filtering on a sequence $\{d_n\}$ to yield $\hat{P}(d_n|f_n, \tau, n\Delta t)$, from which accumulated phase corrupted by measurement noise, $\{y_n\}$, can be obtained from Eq. (1). Since any low pass or a decimation filter has an averaging effect on a signal, decimation filtering a sequence $\{d_n\}$ provides an alternative, software-based approach to physically averaging single shot qubit measurements. Hence, the linear measurement record in Fig. 1(c) arises either from software pre-processing (filtering) data from a single qubit, or from experimental averaging over an ensemble of qubits.

We impose properties on environmental dephasing such that our theoretical designs can enable meaningful predictions. We assume dephasing is non-Markovian, covariance stationary and mean-square ergodic. That is, a single realisation of the process f is drawn from a power spectral density of arbitrary, but non-Markovian form. We further assume that f is a Gaussian process and the separation of timescales between measurement protocols and dephasing dynamics articulated above are met.

Given these conditions, our task is to build a dynamical model to approximately track f over past measurements ($n < 0$), and enable qubit state predictions in future times ($n > 0$). This prediction is represented by the red line in Fig. 1(b-c), and differs from the truth by the so-called estimation (prediction) risk for past (future) times as indicated by shading. We represent our estimate of f for all times using a hat in both the linear and nonlinear measurement models. The major chal-

lenge we face in developing this estimate, \hat{f} (equivalently $\hat{Pr}(d_n|f_n, \tau, n\Delta t)$), is that for a qubit evolving under stochastic dephasing (true state given by black solid line in Fig. 1(b) and (c)), we have no a prior dynamical model for the underlying evolution of f . In the next section, we define the theoretical structure of KF and GPR algorithms which allow us to build that dynamical model directly from the historical measurement record.

III. OVERVIEW OF PREDICTIVE METHODOLOGIES

Our objective is to implement an algorithm permitting learning of underlying qubit dynamics in such a way as to maximize the forward prediction horizon for a given qubit data record. We first quantify the quality of our state estimation procedure. The fidelity of any underlying algorithm during state estimation and prediction, relative to the true state, is expressed by the mathematical quantity known as a Bayes Risk, where zero risk corresponds to perfect estimation. At each time-step, n , the Bayes risk is a mean square distance between truth, f , and prediction, \hat{f} , calculated over an ensemble of M different realisations of true f and noisy data-sets \mathcal{D} :

$$L_{BR}(n|I) \equiv \langle (f_n - \hat{f}_n)^2 \rangle_{f, \mathcal{D}} \quad (2)$$

The notation $L_{BR}(n|I)$ expresses that the Bayes Risk value at n is conditioned on I , a placeholder for free parameters in the design of the predictor, \hat{f}_n . State estimation risk is Bayes Risk incurred during $n \in [-N_T, 0]$; prediction risk is the Bayes Risk incurred during $n \in [0, N_P]$. State estimation and prediction risk regions for one realisation of dephasing noise are shaded in Fig. 1-3. We therefore define the forward prediction horizon as the number of time-steps for $n^* \in [0, N_P]$ during which a predictive algorithm incurs a lower Bayes prediction risk than naively predicting $\hat{f}_n \equiv \mu_f = 0 \quad \forall n$, the mean qubit behaviour under zero-mean dephasing noise.

With this concept in mind, we introduce two general approaches for algorithmic learning relevant to the structures of the problem we have introduced. Our general approach is shared between all algorithms employed and is represented schematically for the KF and GPR in Fig. 2. Stochastic qubit evolution is depicted for one realisation of f [black solid line] given noisy linear measurements [black dots] corrupted by Gaussian white measurement noise v_n . Our overall task is to produce an estimate, given by the red line, which minimizes risk for the prediction period. Ideally both estimation risk and prediction risk are minimized simultaneously for well performing implementations.

Examining the insets in both panels of Fig. 2, both frameworks start with a prior Gaussian distribution over qubit states [purple] that is constrained by the measurement record to yield a posterior Gaussian distribution of the qubit state [red]. The prior captures assumptions

about the qubit state before any data is seen and the posterior expresses our best knowledge of the qubit state under a Bayesian framework. The posterior distribution in both KF and GPR is used to generate qubit state estimates ($n < 0$) and predictions ($n > 0$) [red solid line]. However the computational process by which this posterior is inferred differs significantly between the two methods; we provide an overview of the central features of these algorithms below.

The key feature of a Kalman filter is the recursive learning procedure shown in the inset to Fig. 2(a). Our knowledge of the qubit state is summarised by the prior and a posterior Gaussian probability distributions and these are created and collapsed recursively *at each time step*. The mean of these distributions is the true Kalman state, x_n , and the covariance of these distributions, P_n , captures the uncertainty in our knowledge of x_n ; together both define the Gaussian distribution. The Kalman filter produces an *estimate* of the state, \hat{x}_n at each step through this recursive procedure taking into account two factors. First, the Kalman gain, γ_n , updates our knowledge of (x_n, P_n) within each time step n and serves as a weighting factor for the difference between incoming data, and our best estimate for an observation based on \hat{x}_n , suitably transformed via the measurement action, $h(\hat{x}_n)$. Next, the dynamical model Φ_n propagates the state and covariance, (x_n, P_n) , to the next time step, such that the posterior moments at n define the prior at $n+1$. This process occurs for each time step and an estimate of a true x_n state is built up recursively based on all of our existing knowledge, namely, a linear combination of all past measurements; and all previously generated state estimates. Beyond $n = 0$ we perform predictions in the absence of further measurement data by simply propagating the dynamic model with the Kalman gain set to zero. Full details of the KF algorithm appear below in Section III A.

In our application, we define the Kalman state, x_n , the dynamical model Φ_n , and a measurement action $h(x_n)$ such that the Kalman Filtering framework can track a non-Markovian qubit state trajectory due to an arbitrary realisation of f . In standard KF implementations, the discrete-time sequence $\{x_n\}$, defines a “hidden” signal that cannot be observed, and the dynamic model Φ_n is known. We deviate from this standard construction such that our true Kalman state and its uncertainty, (x_n, P_n) , do not have a direct physical interpretation. Kalman x_n has no a priori deterministic component and corresponds to arbitrary power spectral densities describing f . Hence, the role of the Kalman x_n is to represent an abstract correlated process that, upon measurement, yields physically relevant quantities governing qubit dynamics. Moreover a key challenge described in detail below is to construct an effective Φ_n from the measurement record.

In contrast to the recursive approach taken in the KF, a GPR learning protocol illustrated schematically in Fig. 2(b) selects a *random process* to best describe overall dynamical behaviour of the qubit state under one

propagate the hidden state x_n according to a dynamical model Φ_n corrupted by Gaussian white process noise, w_n .

$$x_n = \Phi_n x_{n-1} + \Gamma_n w_n \quad (3)$$

$$w_n \sim \mathcal{N}(0, \sigma^2) \quad \forall n \quad (4)$$

Process noise has no physical meaning in our application - w_n is shaped by Γ_n and deterministically colored by the dynamical model Φ_n to yield a non-Markovian x_n representing qubit dynamics under generalised environmental dephasing. In addition to coloring via the dynamical model, the process noise covariance matrix, $Q_n \equiv \Gamma_n \Gamma_n^T$, offers an additional mechanism to shape input white noise by designing Γ_n .

We measure x_n using an ideal measurement protocol, $h(x_n)$, and incur additional Gaussian white measurement noise v_n with scalar covariance strength R , yielding scalar noisy observations y_n :

$$y_n = z_n + v_n \quad (5)$$

$$z_n \equiv h(x_n) \quad (6)$$

$$v_n \sim \mathcal{N}(0, R) \quad \forall n \quad (7)$$

The measurement procedure, $h(x_n)$, can be linear or non-linear, allowing us to explore both regimes in our physical application.

With appropriate definitions, the Kalman equations below specify all Kalman algorithms in this paper. At each time step, n , we denote estimates of the moments of the prior and posterior distributions (equivalently, estimates of the true Kalman state) with $(\hat{x}_n(-), \hat{P}_n(-))$ and $(\hat{x}_n(+), \hat{P}_n(+))$ respectively. The Kalman update equations take a generic form (c.f. [31]) :

$$\hat{x}_n(-) = \Phi_{n-1} \hat{x}_{n-1}(+) \quad (8)$$

$$Q_{n-1} = \sigma^2 \Gamma_{n-1} \Gamma_{n-1}^T \quad (9)$$

$$\hat{P}_n(-) = \Phi_{n-1} \hat{P}_{n-1}(+) \Phi_{n-1}^T + Q_{n-1} \quad (10)$$

$$\gamma_n = \hat{P}_n(-) H_n^T (H_n \hat{P}_n(-) H_n^T + R_n)^{-1} \quad (11)$$

$$\hat{y}_n(-) = h(\hat{x}_n(-)) \quad (12)$$

$$\hat{x}_n(+) = \hat{x}_n(-) + \gamma_n (y_n - \hat{y}_n(-)) \quad (13)$$

$$\hat{P}_n(+) = [1 - \gamma_n H_n] \hat{P}_n(-) \quad (14)$$

To reiterate, Eq. (8) and Eq. (10) bring the best state of knowledge from the previous time step into the current time step, n , as a prior distribution. Dynamical evolution is modified by features of process noise, as encoded in Eq. (9), and propagated in Eq. (10). The propagation of the moments of the a priori distribution, as outlined thus far, does not depend on the incoming measurement, y_n , but is determined entirely by the a priori (known) dynamical model, in our case $\Phi \equiv \Phi_n, \forall n$.

The Kalman gain in Eq. (11) depends on the uncertainty in the true state, $\hat{P}_n(-)$ and is modified by features of the measurement model, H_n , and measurement noise, $R_n \equiv R, \forall n$. It serves as an effective weighting

function for each incoming observation. Before seeing any new measurement data, the filter predicts an observation $\hat{y}_n(-)$ corresponding to the best available knowledge at n in Eq. (12). This value is compared to the actual noisy measurement y_n received at n , and the difference is used to update our knowledge of the true state via Eq. (13). If measurement data is noisy and unreliable (high R), then γ has a small value, and the algorithm propagates Kalman state estimates according to the dynamical model and effectively ignores data. In particular, only the second terms in both Eq. (13) and Eq. (14) represent the Bayesian update of the moments of a prior distribution ((-) terms) to the posterior distribution ((+) terms) at n . If $\gamma_n \equiv 0$, then the prior and posterior moments at any time step are exactly identical by Eqs. (13) and (14), and only dynamical evolution occurs using Eqs. (8) to (10). This is the condition we employ when we seek to make forward predictions beyond a single time-step, and hence we set $\gamma \equiv 0$ during future prediction.

Since we do not have a known dynamical model Φ for describing stochastic qubit dynamics under f , we will need to make design choices for $\{x, \Phi, h(x), \Gamma\}$ such that f can be approximately tracked. These design choices will completely specify algorithms introduced below and represent key findings with respect to our work in this manuscript. For a linear measurement record, $h(x) \mapsto Hx$ and we compare predictive performance for Φ modeling stochastic dynamics either via so-called ‘autoregressive’ processes in the AKF, or via projection onto a collection of oscillators in the LKKFB. In addition, we use the dynamics of AKF to define a Quantised Kalman filter (QKF) with a non-linear, quantised measurement model such that the filter can act directly on binary qubit outcomes. We provide the relevant details in sub-sections below.

1. Autoregressive Kalman Filter (AKF)

Recursive autoregressive methods are well-studied in classical control applications (c.f. [32]) presenting opportunities to leverage existing engineering knowledge in developing quantum control strategies. In our application, we use an autoregressive Kalman filter to probe arbitrary, covariance-stationary qubit dynamics such that the dynamic model is constructed as a weighted sum of q past values driven by white noise *i.e.* an autoregressive process of order q , AR(q). Using Wold’s decomposition, it can be shown that any zero mean covariance stationary process representing qubit dynamics has a representation in the mean-square limit by an autoregressive process of finite order, as in Appendix B.

The study of AR(q) processes falls under a general class of techniques based on autoregressive moving average (ARMA) models in adaptive control engineering and econometrics (e.g. [33, 34] respectively). For high- q models in a typical time-series analysis, it is possible to

decompose an $\text{AR}(q)$ into an ARMA model with a small number of parameters [35, 36]. However, we retain a high- q model to probe arbitrary power spectral densities. Further, literature suggests employing a high- q model is relatively easier than a full ARMA estimation problem and enables lower prediction errors [35, 37].

To construct the Kalman dynamical operator Φ for the AKF, we introduce a set of q coefficients $\{\phi_{q' \leq q}\}$, $q' = 1, \dots, q$ to specify the dynamical model:

$$f_n = \phi_1 f_{n-1} + \phi_2 f_{n-2} + \dots + \phi_q f_{n-q} + w_n \quad (15)$$

We thus see that the dynamical model is constructed as a weighted sum of time-retarded samples of f , with weighting factors given by the autoregressive coefficients up to order (and hence time lag) q . For small $q < 3$, it is possible to extract simple conditions on the coefficients, $\{\phi_{q' \leq q}\}$, that guarantee properties of f : for example, that f is covariance stationary and mean square ergodic. In our application, we freely employ arbitrary- q models via machine learning in order to improve our approximation of an arbitrary f . Any $\text{AR}(q)$ process can be recast (non-uniquely) into state space form ([4]), and we define the AKF by the following substitutions into Kalman equations:

$$x_n \equiv [f_n \dots f_{n-q+1}]^T \quad (16)$$

$$\Gamma_n w_n \equiv [w_n 0 \dots 0]^T \quad (17)$$

$$\Phi_{AKF} \equiv \begin{bmatrix} \phi_1 & \phi_2 & \dots & \phi_{q-1} & \phi_q \\ 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \ddots & \vdots & \vdots \\ 0 & 0 & \ddots & 0 & 0 \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix} \quad \forall n \quad (18)$$

$$H \equiv [1 \ 0 \ 0 \ 0 \dots 0] \quad \forall n \quad (19)$$

The matrix Φ_{AKF} is the dynamical model used to recursively propagate the unknown state during state estimation in the AKF, as represented schematically in the upper half of Fig. 3. In general, the $\{\phi_{q' \leq q}\}$ employed in Φ_{AKF} must be learned through an optimisation procedure using the measurement record, where the set of parameters to be optimised is $\{\phi_1, \dots, \phi_q, \sigma^2, R\}$. This procedure yields the optimal configuration of the autoregressive Kalman filter, but at the computational cost of a $q + 2$ -dimensional Bayesian learning problem for arbitrarily large q .

The Least Squares Filter (LSF) in [28] considers a weighted sum of past measurements to predict the i -th step ahead measurement outcome, $i \in [0, N_P]$. A gradient descent algorithm learns the weights, $\{\phi_{q' \leq q}\}$ for the previous q past measurements, and a constant offset value for non-zero mean processes, to calculate the i -th step ahead prediction. The set of N_P LSF models, collectively, define the set of predicted qubit states under an LSF acting on a measurement record. For $i = 1$, equivalent to the single-step update employed in the Kalman filter, we assert that learned $\{\phi_{q' \leq q}\}$ in LSF effectively

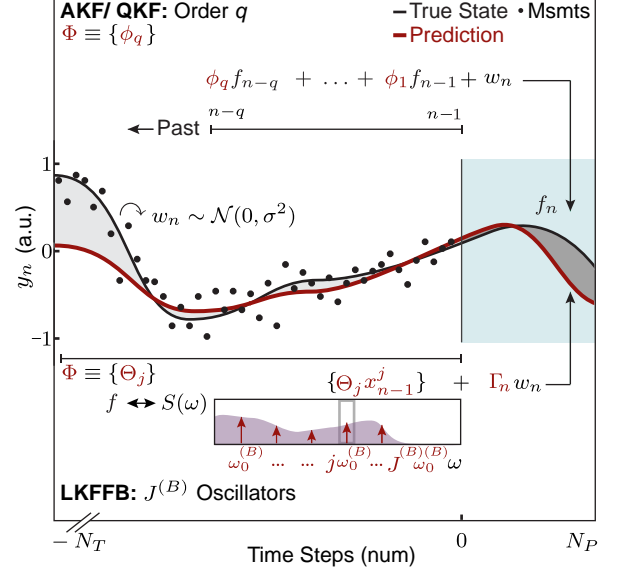


FIG. 3. Approaches to construction of the KF dynamical model. Panel (a) from Fig. 2 superimposed with Kalman dynamical models, $\Phi \equiv \Phi_n, \forall n$. (a) AKF/QKF: A set of autoregressive coefficients, $\{\phi_{q' \leq q}\}$, define Φ to yield f_n as a weight sum of q past measurements. (b) LKFFB: Red arrows with heights $\|x_n^j\|$ depict set of basis oscillators for $j = 1, \dots, J^{(B)}$ probe true purple spectrum of f_n and yields time domain dynamics of f_n as a stacked system of resonators, Θ_j . Black L-shaped arrows depict a single instance of f_n at $n = 0$ based on historical $\{f_{n-1}, f_{n-2}, \dots\}$.

implements an $\text{AR}(q)$ process (we validate numerically in Section IV). Under this condition, and for zero-mean w_n , the LSF in [28] by definition searches for coefficients for the weighted linear sum of past q measurements, as described in in Eq. (15).

We use the parameters $\{\phi_{q' \leq q}\}$ learned in the LSF to define Φ in Eq. (18), therefore reducing the computational complexity of the remaining optimisation from $((q + 2) \rightarrow 2)$ -dimensional for an AKF of order q . Since Kalman noise parameters (σ^2, R) are subsequently auto-tuned using a Bayes Risk optimisation procedure (see Section IV A), we optimise over potential remaining model errors and measurement noise.

In general, LSF performance improves as q increases and a full characterisation of model-selection decisions for LSF are given in [28]. Defining an absolute value for the optimal q is somewhat arbitrary as it is defined relative to the extent to which a true f is oversampled in the measurement routine and the finite size of the data. For all analyses presented here, we fix the ratio $q\Delta t = 0.1$ (a.u.) and $q/N_T = 0.05$ (a.u.), where the experimental sampling rate is $1/\Delta t$, N_T and $\{\phi_{q' \leq q}\}$ are identical in the AKF and LSF. In practice this ensures numerical convergence of the LSF during training.

2. Liska Kalman Filter with Fixed Basis (LKFFB)

In LKFFB, we effectively perform a Fourier decomposition of the underlying f in order to build the dynamic model, Φ , for the Kalman filter. Here, we project our measurement record on $J^{(B)}$ oscillators with fixed frequency $\omega_j \equiv j\omega_0^{(B)}$ with j an integer as $j = 1, \dots, J^{(B)}$. The temporal resolution of the state tracking procedure is set by the maximum frequency in the selected basis and properties of the spacing between adjacent basis frequencies. The superscript $^{(B)}$ indicates Fourier domain information about an algorithmic basis, as opposed to information about the true (unknown) dephasing process. The LKFFB allows instantaneous amplitude and phase tracking for each basis oscillator, directly enabling forward prediction from the learned dynamics. The structure of this Kalman filter, referred to as the Liska Kalman Filter (LKF), was developed in [29]; adding a fixed basis in this application yields the Liska Kalman Filter with a Fixed Basis (LKFFB).

For our application, the true hidden Kalman state, x , is encoded as a collection of sub-states, x^j , for the j^{th} oscillator. For clarity we remind that the superscript is used as an index rather than a power. Each sub-state is labeled by a real and imaginary component which we represent in vector notation:

$$x_n \equiv [x_n^1 \dots x_n^j \dots x_n^{J^{(B)}}] \quad (20)$$

$$A_n^j \equiv \text{Re}(x_n^j) \quad (21)$$

$$B_n^j \equiv \text{Im}(x_n^j) \quad (22)$$

$$x_n^j \equiv \begin{bmatrix} A_n^j \\ B_n^j \end{bmatrix} \quad (23)$$

The algorithm tracks the real and imaginary parts of the Kalman sub-state simultaneously in order calculate the instantaneous amplitudes ($\|x_n^j\|$) and phases (θ_n^j) for each Fourier component:

$$\|x_n^j\| \equiv \sqrt{(A_n^j)^2 + (B_n^j)^2} \quad (24)$$

$$\theta_n^j \equiv \tan^{-1} \frac{B_n^j}{A_n^j} \quad (25)$$

The dynamical model for LKFFB is now constructed as a stacked collection of these independent oscillators. The sub-state dynamics match the formalism of a Markovian stochastic process defined on a circle for each basis frequency, ω_j , as in Ref. [38]. We stack $\Theta(j\omega_0^{(B)}\Delta t)$ for all ω_j along the diagonal to obtain the full dynamical matrix for Φ_n :

$$\Phi_n \equiv \begin{bmatrix} \Theta(\omega_0^{(B)}\Delta t) & \dots & 0 \\ \dots & \Theta(j\omega_0^{(B)}\Delta t) & \dots \\ 0 & \dots & \Theta(J^{(B)}\omega_0^{(B)}\Delta t) \end{bmatrix} \quad (26)$$

$$\Theta(j\omega_0^{(B)}\Delta t) \equiv \begin{bmatrix} \cos(j\omega_0^{(B)}\Delta t) & -\sin(j\omega_0^{(B)}\Delta t) \\ \sin(j\omega_0^{(B)}\Delta t) & \cos(j\omega_0^{(B)}\Delta t) \end{bmatrix} \quad (27)$$

We obtain a single estimate of the true hidden state by defining the measurement model, H , by concatenating $J^{(B)}$ copies of the row vector $[1 \ 0]$:

$$H \equiv [1 \ 0 \dots 1 \ 0 \dots 1 \ 0] \quad (28)$$

Here, the unity values of H pick out and sum the Kalman estimate for the real components of f while ignoring the imaginary components, namely, we sum A_n^j for all $J^{(B)}$ basis oscillators.

In [29], a state-dependent process-noise-shaping matrix is introduced to enable potentially non-stationary instantaneous amplitude tracking in LKFFB for each individual oscillator:

$$\Gamma_{n-1} \equiv \Phi_{n-1} \frac{x_{n-1}}{\|x_{n-1}\|} \quad (29)$$

For the scope of this manuscript, we retain the form of Γ_n in our application even if true qubit dynamics are covariance stationary. As such, Γ_n depends on the state estimates x . For this choice of Γ_n , we deviate from classical Kalman filters because recursive equations for P cannot be propagated in the absence of measurement data. Consequently, Kalman gains cannot be pre-computed prior to experimental data collection. Details of gain pre-computation in classical Kalman filtering can be found in standard textbooks (e.g. [31]).

There are two ways to conduct forward prediction for LKFFB and both are numerically equivalent for an appropriate choice of basis: (i) we set the Kalman gain to zero and recursively propagate using Φ ; (ii) we define a harmonic sum using the basis frequencies and learned $\{\|x_n^j\|, \theta_n^j\}$. This harmonic sum can be evaluated for all future time to yield forward predictions in a single calculation. The choice of basis for an LKFFB and its implications for optimal predictive performance are discussed in Appendix C 2.

3. Quantised Kalman Filter (QKF)

In QKF, we implement a Kalman filter that acts directly on discretised measurement outcomes, $d \in \{0, 1\}$. To reiterate the discussion of Fig. 1(a), this means that the measurement action in QKF must be non-linear and take as input quantised measurement data. This holds true irrespective of our dynamical model, Φ . In our application we set the dynamical model to be identical to that employed in the AKF, allowing isolation of the effect of the nonlinear, quantised measurement action.

With unified notation across AKF and QKF, we define a non-linear measurement model $h(x)$ and its Jacobian, H as:

$$z_n \equiv h(x_n[0]) \equiv \frac{1}{2} \cos(f_n) \quad (30)$$

$$\Rightarrow H_n \equiv \frac{dh(f_n)}{df_n} = -\frac{1}{2} \sin(f_n) \quad (31)$$

During filtering, $z_n = h(x_n[0])$ is used to compute measurement residuals when updating the true Kalman state, x_n , whereas the state variance estimate, P_n , is propagated using the Jacobian, H_n . Further, the Jacobian is used to compute the Kalman gain. Hence the filter can quickly destabilise if the linearisation of $h(\cdot)$ by H_n doesn't hold during dynamical propagation, resulting in a rapid build up of errors.

In this construction, the entity z_n is associated with an abstract 'signal': a sequence formed by repeated applications of the likelihood function for a single qubit measurements in Eq. (1). The true stochastic qubit phase, f_n , is our Kalman hidden state, x_n . Subsequently, we extract an estimate of the true bias, z_n , as an unnatural association of the Kalman measurement model with Born's rule. The sequence $\{z_n\}$ is not observable, but can only be inferred over a large number of experimental runs.

To complete the measurement action, we implement a biased coin flip within the QKF filter given \tilde{y}_n . While the qubit provides measurement outcomes which are naturally quantised, we require a theoretical model, \mathcal{Q} , to generate quantised measurement outcomes with statistics that are consistent with Born's rule in order to propagate the dynamic Kalman filtering equations appropriately. In order to build this machinery we modify the procedure in [39] to quantise z_n using biased coin flips. In our notation, we represent a black-box quantiser, \mathcal{Q} , that gives only a 0 or a 1 outcome based on \tilde{y}_n :

$$d_n = \mathcal{Q}(\tilde{y}_n) \quad (32)$$

$$= \mathcal{Q}(h(f_n) + v_n) \quad (33)$$

The use of the notation \tilde{y}_n is meant to indicate a correspondence with y_n introduced earlier, while the physical meaning differs due to the discretised nature of the QKF. Therefore, the stochastic changes in $\{\tilde{y}_n\}$ are represented in the bias of a coin flip, subject to proper normalisation constraints which maintains $|\tilde{y}_n| \leq 0.5$:

$$Pr(d_n|\tilde{y}_n, f_n, \tau) \equiv \mathcal{B}(n_B = 1; p_B = \tilde{y}_n + 0.5) \quad (34)$$

QKF uses Eq. (34) to define a biased coin-flip during filtering, where n_B represents a single coin flip, p_B represents the stochastically drifting bias on the coin. Kalman filtering with the coin-flip quantisation defined by Eq. (34) presents a departure from classical amplitude quantisation procedures in [39, 40].

From a computational perspective, we modify the process noise features definition from AKF to QKF. We set $Q \equiv \sigma^2 \Gamma \Gamma^T \rightarrow \sigma^2 \mathcal{I} \quad \forall n, \mathcal{I}$ is $q \times q$ identity matrix, from AKF to QKF. This rationale for this modification is that it smears out the effect of white process noise in a way that stabilizes inversions in the gain calculation in the Kalman filter, but does not correlate any two Kalman states in time (diagonal matrix). In practice, this modification only yields mild improvements over the original AKF process noise features matrix.

The definitions of $\{\mathcal{Q}, h(x_n), H_n, Q\}$ in this subsection, and dynamics $\{x_n, \Phi\}$ from the AKF now completely

specify the QKF algorithm for application to a discrete, single-shot measurement record as depicted in Fig. 1 (a).

B. Gaussian Process Regression (GPR)

In GPR, correlations in the measurement record can be learned if one projects data on a distribution of Gaussian processes, $Pr(f)$ with an appropriate encoding of their covariance relations via a kernel, $\Sigma_f^{n_1, n_2}$. We return to the linear measurement record and the definition of scalar noisy observations y_n corrupted by Gaussian measurement noise, v_n , as considered previously for AKF, LSF, and LKFFB. Under linear operations, the distribution of measured outcomes, y_n , is also a Gaussian. The mean and variance of $Pr(y)$ depends on the mean μ_f and variance Σ_f of the prior $Pr(f)$, and the mean $\mu_v \equiv 0$ and variance R of the measurement noise:

$$f \sim Pr_f(\mu_f, \Sigma_f) \quad (35)$$

$$y \sim Pr_y(\mu_f, \Sigma_f + R) \quad (36)$$

For covariance stationary f , correlation relationships depend solely on the time lag, $\nu \equiv \Delta t|n_1 - n_2|$ between any two time points $n_1, n_2 \in [-N_T, N_P]$. An element of the covariance matrix, $\Sigma_f^{n_1, n_2}$, corresponds to one value of lag, ν , and the correlation for any given ν is specified by the covariance function, $R(\nu)$:

$$\Sigma_f^{n_1, n_2} \equiv R(\nu) \quad (37)$$

Any unknown parameters in the encoding of correlation relations via $R(\nu)$ are learned by solving the optimisation problem outlined in Section IV A. The optimised GPR model is then applied to datasets corresponding to new realisations of f . Let indices $n \in N_T \equiv [-N_T, 0]$ denote training points, and let a length N^\dagger vector contain arbitrary testing points $n^\dagger \in [-N_T, N_P]$. These testing points in machine learning language encompass both state estimation and prediction points in our notation. We now define the joint distribution $Pr(y, f^\dagger)$, where f^\dagger represents the true process evaluated by GPR at desired test points:

$$\begin{bmatrix} f^\dagger \\ y \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mu_{f^\dagger} \\ \mu_y \end{bmatrix}, \begin{bmatrix} K(N^\dagger, N^\dagger) & K(N_T, N^\dagger) \\ K(N^\dagger, N_T) & K(N_T, N_T) + R \end{bmatrix}\right) \quad (38)$$

The additional 'kernel' notation $\Sigma_f \equiv K(N_T, N_T)$ is ubiquitous in GPR. Time domain correlations specified by $R(\nu)$ populate each element of a matrix $K(\cdot, \cdot)$, where the dimensions of the matrix depend on the vector length of each argument. For example, for $K(N_T, N_T)$, the notation defines a square matrix where diagonals correspond to $\nu = 0$ and off-diagonal elements correspond to separation of two arbitrary points in time i.e. $\nu \neq 0$.

Following [41], the moments of the conditional predictive distribution $Pr(f^\dagger|y)$ can be derived from the joint

distribution $Pr(y, f^\dagger)$ via standard Gaussian identities:

$$\mu_{f^\dagger|y} = \mu_f + K(N^\dagger, N_T)(K(N_T, N_T) + R)^{-1}(y - \mu_y) \quad (39)$$

$$\Sigma_{f^\dagger|y} = K(N^\dagger, N^\dagger) - K(N^\dagger, N_T)(K(N_T, N_T) + R)^{-1}K(N_T, N^\dagger) \quad (40)$$

The prediction procedure outlined above holds true for any choice of kernel, $R(\nu)$. In any GPR implementation, the dataset, y , constrains the prior model yielding an a posteriori predictive distribution. The mean values of this predictive distribution, $\mu_{f^\dagger|y}$, are the state predictions for the qubit under dephasing at test points in N^\dagger .

In our work we focus on a ‘periodic kernel’ to encode a covariance function which is theoretically guaranteed to approximate any zero-mean covariance stationary process, f , in the mean square limit, by having the same structure as a covariance function for trigonometric polynomials with infinite harmonic terms [38, 42]. The sine squared exponential kernel represents an infinite basis of oscillators and is defined as:

$$R(\nu) \equiv \sigma^2 \exp\left(-\frac{2 \sin^2\left(\frac{\omega_0^{(B)} \nu}{2}\right)}{l^2}\right) \quad (41)$$

This kernel is described using just two key hyperparameters: the frequency-comb spacing for our infinite basis of oscillators, ω_0 , and a dimensionless length scale, l . We use physical sampling considerations to approximate their initial conditions prior to an optimisation procedure, namely, that the longest correlation length encoded in the data sets the frequency resolution of the comb, and the scale at which changes in f are resolved is limited physically by the minimum time taken between sequential Ramsey measurements:

$$\frac{\omega_0^{(B)}}{2\pi} \sim \frac{1}{\Delta t N} \quad (42)$$

$$l \sim \Delta t \quad (43)$$

Because the periodic kernel can be shown to be formally equivalent to the basis of oscillators employed in the LKFFB algorithm in a limiting case (see Appendix C for a discussion using results in [42]), the inclusion of GPR using this kernel permits a comparison of the underlying algorithmic structures for the task of predictive estimation using spectral methods.

For the analysis of covariance stationary time series under a GPR framework, we de-emphasise popular kernel choices such as: a Gaussian kernel (RBF), a scale mixture of Gaussian kernels (RQ), and Matern kernels (e.g. MAT32) [41, 43]. An arbitrary-scale mixture of zero-mean Gaussian kernels will probe an arbitrary area around zero in the Fourier domain, as schematically depicted in Fig. 2(a). While such kernels capture the continuity assumption ubiquitous in machine learning, they

are structurally inappropriate for probing a process characterized by an arbitrary power spectral density (e.g. ohmic noise). Another common kernel for time-series analysis is a quasi periodic kernel (QPER) defined by a product of an RBF with a periodic kernel [44]. This corresponds to a convolution in the Fourier domain giving rise to a comb of Gaussians at the expense of an increase in the number of parameters required for kernel tuning. One can also consider specific types of $AR(q)$ processes using Matern kernels of order $q + 1/2$ but with increased restrictions on the form of coefficients [41, 45]. A simple consideration of autoregressive approaches suggest that a Matern kernel for $q = 1$ (MAT32) can be briefly trialed under GPR, whereas high- q autoregressive processes are naturally and generally treated under a KF framework. Further discussion of kernel choice appears in Sec. V.

IV. ALGORITHM PERFORMANCE CHARACTERISATION

In the results to follow, our metric for characterising performance of optimally tuned algorithms will be the normalised Bayes prediction risk:

$$\tilde{L}_{BR} \equiv \frac{L_{BR}(n|I)}{\langle (f_n - \mu_f)^2 \rangle_{f,D}}, \quad \mu_f \equiv 0 \quad (44)$$

A desirable forward prediction horizon corresponds to maximal $n^* \in [0, N_P]$ for which normalised Bayes prediction risk at all time-steps $n \leq n^*$ is less than unity. We compare the difference in maximal forward prediction horizons between algorithms in the context of realistic operating scenarios. We begin here by introducing the numerical methods employed for generating data-sets on which predictive estimation is performed.

We simulate environmental dephasing through a Fourier-domain procedure described in Appendix A 2 [46] in order to simulate an f which is mean-square ergodic and covariance stationary. For the results in this manuscript, we choose a flat top spectrum with a sharp high-frequency cutoff for simplicity as this choice of a power spectral density theoretically favors no particular choice of algorithm but violates the Markov property.

In our simulations we also must mimic a measurement process which samples the underlying ‘true’ dephasing process. The algorithmic parameters $\{N_T, \Delta t\}$ represent a sampling rate and Fourier resolution set by the simulated measurement protocol; we choose regimes where the Nyquist rate, $r \gg 2$. In generating noisy simulated measurement records, we corrupt a noiseless measurement by additive Gaussian white noise. Since f is Gaussian, the measurement noise level, $N.L.$, is defined as a ratio between the standard deviation of additive Gaussian measurement noise, \sqrt{R} and the maximal spread of random variables in any realisation f . We approximate the maximal spread of f as three sample standard deviations of one realisation of true f , $N.L. = \sqrt{R}/3\sqrt{\hat{\Sigma}_f^{n,n}}$.

The use of a hat in this notation denotes sample statistics. This computational procedure enables a consistent application of measurement noise for f from arbitrary, non-Markovian power spectral densities. For the case where binary outcomes are required, we apply a biased coin flip using Eq. (34).

A. Algorithmic Optimisation

All algorithms in this manuscript employ machine learning principles to tune unknown design parameters based on training data-sets. The physical intuition associated with optimising our filters is that we are cycling through a large class of general models for environmental dephasing and seeking the model(s) which best fit the data subject to various constraints. This allows each filter to track stochastic qubit dynamics under arbitrary covariance-stationary, non-Markovian dephasing. We elected to deploy an optimisation routine with minimal computational complexity to enable nimble deployment of KF and GPR algorithms in realistic laboratory settings, particularly since LSF optimisation is extremely rapid for our application [28].

Kalman filtering in our setting poses a significant challenge for general optimisers as the lack of theoretical bounds on the values of (σ, R) result in large, flat regions of the Bayes Risk function. Further, the recursive structure of the Kalman filter means that no analytical gradients are accessible for optimising a choice of cost function and a large computational burden is incurred for any optimisation procedure. We randomly distribute (σ_k, R_k) pairs for $k = 1, \dots, K$ over ten orders of magnitude in two dimensions in order to sample the optimisation space.

We then generate a sequence of loss values $L(\sigma_k, R_k)$ for each k by considering a small region around $n = 0$, where the size of the region is n_L number of time steps we look forward or backwards from $n = 0$:

$$L(\sigma_k, R_k) \equiv \sum_{n=1}^{n_L} L_{BR}(n|I = \{\sigma_k, R_k\}). \quad (45)$$

Here, $L_{BR}(n|I = \{\sigma_k, R_k\})$ is given by Eq. (2) and it is summed over $0 \leq n_L \leq |N_T|$ ($0 \leq n_L \leq |N_P|$) backwards (forwards) time-steps for state estimation (prediction). In the notation for I above, we omit Kalman dynamical model design parameters for ease of reading. Typically I would include, for instance, the set of autoregressive coefficients in AKF and the set of fixed basis frequencies in LKFFB. Values of n_L are chosen such that the sequence $\{L(\sigma_k, R_k)\}$ defines sensible shapes of the total loss function over parameter space and the numerical experiments in this manuscript. A choice of small n_L in state-estimation ensures that data near the prediction horizon are employed - a region where the Kalman filter is most likely to have converged. Similarly, in state prediction, large n_L will flatten the true prediction loss function as long-term prediction errors dominate smaller

loss values occurring during the short term prediction period. In addition, one can weight state estimation and state prediction loss functions differently by choosing different values of n_L for state estimation and prediction, though we set n_L to be the same in both regions. While simple and by no means optimal, our tuning approach is computationally tractable and efficient compared to the application of standard optimisation routines where each loss value calculation requires a recursive filter to act on a long measurement record. Further, our approach ensures tuning procedures are performed off-line such that a tuned algorithm is simple in its recursive structure and performs rapid calculations at each time-step.

An ideal parameter pair (σ^*, R^*) minimises Bayes risk over K trials for both state estimation and prediction. We define acceptable low loss regions for state estimation and prediction as being the set which returns loss less than 10% of the median risk over K trials. In the event that low risk regions do not exist for both state estimation and prediction for a given parameter pair, we deem the optimisation to have failed as state estimation performance is uncorrelated with forward prediction (for illustration, see panel (h) of Fig. 7).

In GPR the set of parameters $I = \{\sigma, R, \omega_0^{(B)}, l\}$ requires optimisation. However, in contrast to the KF, no recursion exists and analytic gradients are accessible to simplify the overall optimisation problem. Instead of minimising Bayes state-estimation risk, we follow a popular practice of maximising the Bayesian likelihood. Initial conditions and optimisation constraints are derived from physical arguments as described in Section III.

B. Performance of the KF using linear measurement

The general performance of the various KF algorithms discussed above is illustrated in Fig. 4 which compares the AKF and LKFFB algorithms using a linear measurement record. Here the solid black line represents the underlying true f and solid markers indicate noisy simulated linear measurement data. Future predictions using the various KF formalisms and the (non-recursive) LSF filter [28] are shown as coloured open markers, based on these data. The selected single realization of the prediction process demonstrated in (a) is representative of a broad ensemble of simulated data sets and demonstrates the ability of all algorithms to perform future prediction with varying degrees of success.

In general, our objective is to maximise the forward prediction horizon, n^* , in any algorithmic setting. In Fig. 4(b)-(d), we explore the key determining factors setting the value of the prediction horizon under the three main Kalman filtering algorithms treated here. We plot the ensemble-averaged \bar{L}_{BR} as a function of forward prediction time when adjusting the ratio of the cutoff frequency in the noise, $J\omega_0$, to the sample rate in the measurement routine ($\omega_{(S)} = 2\pi/\Delta t$) without physical

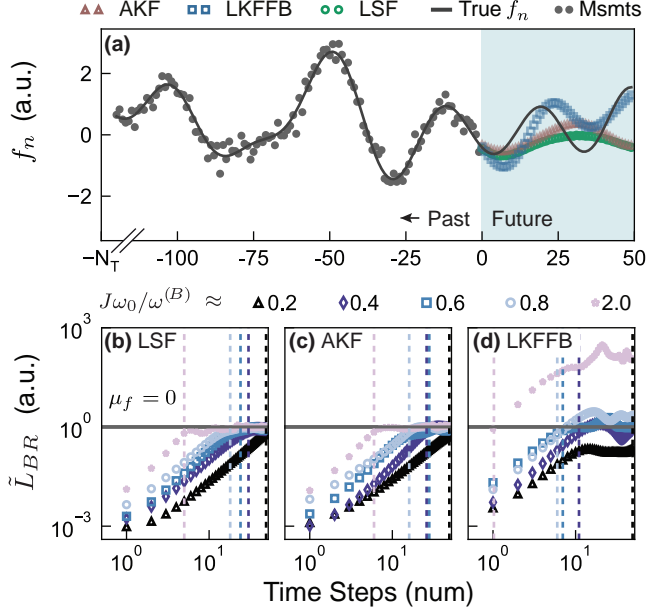


FIG. 4. (a) Solid dots depict y_n against time-steps n and data collection ceases at $n = 0$. Optimised LSF, AKF and LKFFB yield predictions $n > 0$ in the blue region plotted as open, coloured markers. A black solid line shows one realisation of true f_n , drawn from a flat top spectrum with J true Fourier components spaced ω_0 apart and uniformly randomised phases. Other parameters: $\omega_0/\omega_0^{(B)} \notin \mathbb{Z}$ (natural numbers), $J = 45000$, $\omega_0/2\pi = \frac{8}{9} \times 10^{-3}$ Hz such that > 500 number of true components fall between adjacent LKFFB oscillators; $N.L. = 10\%$. (b)-(d) Procedure in (a) is repeated for ensemble M different realisations of f and noisy datasets to compute \tilde{L}_{BR} for LSF, AKF, and LKFFB. \tilde{L}_{BR} v. $n \in [0, N_P]$ is plotted; dark-grey horizontal line marks $\tilde{L}_{BR} \equiv 1$ for predicting the mean $\mu_f \equiv 0$. Vertical dashed lines mark the forward prediction horizon, n^* , where $\tilde{L}_{BR} \lesssim 0.8 < 1$ for all prediction time steps $0 < n \leq n^*$ in out-performing predicting the noise mean. Marker color (dark indigo to pink) depicts true f cutoff, $J\omega_0$ varied relative to $\omega^{(B)} \equiv \omega_0^{(B)} J^{(B)} \approx r\omega_{(S)}$, with fixed Nyquist $r \gg 2$; $\omega_0/2\pi = 0.497$ Hz, $J = 20, 40, 60, 80, 200$; $N.L. = 1\%$. For all (a)-(d), a trained LKFFB is implemented with $\omega_0^{(B)}/2\pi = 0.5$ Hz and $J^{(B)} = 100$ oscillators; trained AKF / LSF models are $q = 100$; with $N_T = 2000$, $N_P = 50$ steps, $\Delta t = 0.001$ s, $M = 50$ runs, $K = 75$ optimisation trials.

aliasing such that Nyquist $r \gg 2$ and $\omega_{(S)} \approx \omega^{(B)}/r$, where $\omega^{(B)}$ incorporates a (potentially incorrect) bandwidth assumption about dephasing noise for LKFFB. Here again, we have a forward prediction horizon for time-steps $0 < n < n^*$ if $\tilde{L}_{BR} \lesssim 1$ for all time-steps in this region and an algorithm seeks to maximise n^* . In this region, each algorithm predicts future dynamics better than naively predicting the mean behaviour of f ($\mu_f \equiv 0$), indicated by a dark-grey horizontal line.

The prediction horizon, indicated approximately by dashed vertical lines, for all algorithms increases as the measurement becomes sufficiently fast to sample the

highest frequency dynamics of f . We confirm numerically that absolute prediction horizons for any algorithm are arbitrary and adjustable through the sample rate, allowing us to restrict our analysis to comparative statements between algorithms for future results. While differences between protocols appear reasonably small we note that in most cases examined the AKF demonstrates superior performance to the LKFFB subject to the realistic constraint that the true dynamics of f cannot be perfectly projected onto the basis used in LKFFB (the latter situation corresponding to substantial a priori knowledge of the dynamics of f). The role of undersampling in the LKFFB becomes pronounced as predictive estimates lead to unstable behavior relative to the naive prediction of $\mu_f = 0$ in the case $J\omega_0/\omega^{(B)} = 2$ in Fig. 4(d). The AKF and LSF share autoregressive coefficients and therefore both algorithms demonstrate comparable \tilde{L}_{BR} prediction risk in the ensemble average.

A key implied benefit of the use of Kalman filtering vs the LSF with high-order autoregressive dynamics alone is the addition of robustness against measure-

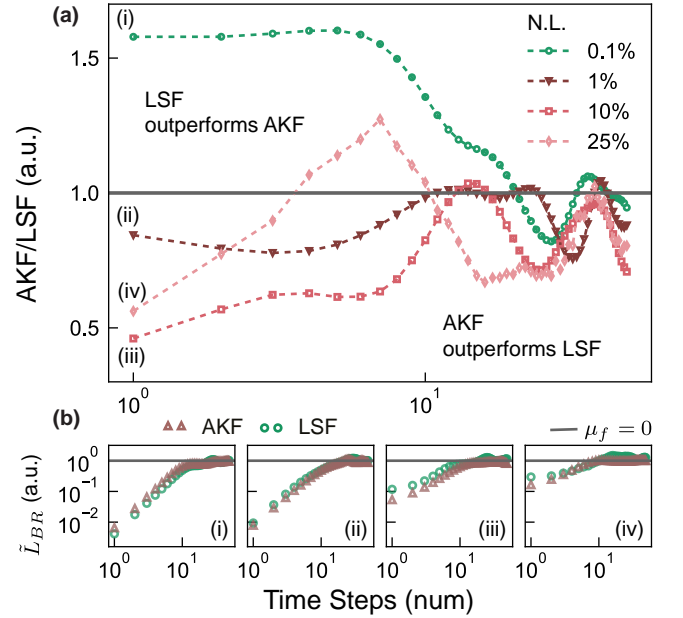


FIG. 5. Measurement noise filtering in AKF v. LSF. (a) Dashed-lines with markers depict the ratio of \tilde{L}_{BR} for AKF to LSF against time-steps $n > 0$; for cases (i)-(iv) with $N.L. = 0.1, 1.0, 10.0, 25.0\%$. Green trajectory shows LSF outperforms AKF with ratio > 1 for $n \leq n^*$; crimson trajectories show AKF outperforms LSF with ratio < 1 for $n \leq n^*$. (b) \tilde{L}_{BR} against n is plotted for cases (i)-(iv) confirms a maximal forward prediction horizon marked by n^* , exists for all ratios in (a) for both LSF and AKF. In (a) and (b), AKF and LSF share identical $\{\phi_q\}$. True f is drawn from a flat top spectrum with $\omega_0/2\pi = \frac{8}{9} \times 10^{-3}$ Hz, $J = 45000$, $N_T = 2000$, $N_P = 100$ steps, $\Delta t = 0.001$ s, $r = 20$ such that Fig. 6(c) corresponds to case (ii) in this figure. AKF is optimised with $q = 100$, $M = 50$ runs, $K = 75$ trials.

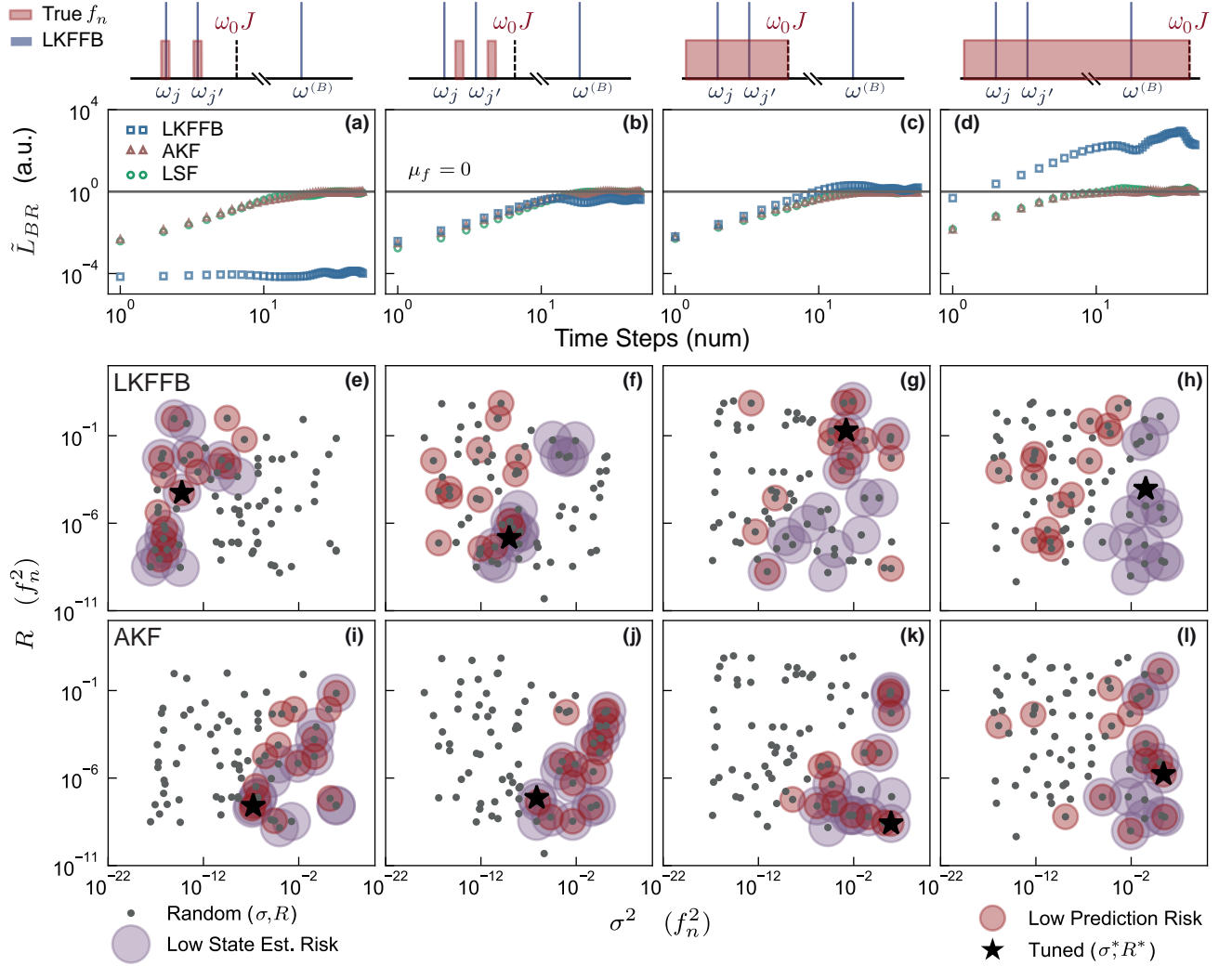


FIG. 6. Comparison of KF performance under various imperfect learning scenarios. (a)-(d) True noise properties are varied to introduce pathological learning with respect to fixed algorithmic configuration: $\omega_0/2\pi = 0.5, 0.499, \frac{8}{9} \times 10^{-3}, \frac{8}{9} \times 10^{-3}$ Hz and $J = 80, 80, 45000, 80000$ respectively. The relationship between LKFFB basis and true noise spectrum is shown schematically above columns: (a) perfect learning; (b) imperfect projection on LKFFB basis; (c) finite computational Fourier resolution; (d) relaxed basis bandwidth assumption. (a)-(d) \tilde{L}_{BR} against time-steps $n > 0$ is shown for LKFFB, AKF, and LSF. (e)-(l) Optimisation results for LKFFB [top row] and AKF [bottom row] in each of the four regimes in (a)-(d). Grey dots depict K random (σ^2, R) pairs; where M realisations of f, \mathcal{D} are used to calculate \tilde{L}_{BR} for each pair. Purple (crimson) circles represent low loss regions where risk value in Eq. (45), for (σ^2, R) is $< 10\%$ of the median risk value during state estimation (prediction) for $-n_L < n < 0$ ($n_L > n > 0$), with $n_L = 50$. Black star, (σ^*, R^*) , minimises risk values over purple circles during state estimation. A KF filter is ‘tuned’ if optimal (σ^*, R^*) lies in the overlap of low loss regions for state estimation [purple] and prediction [crimson]; disjoint regions in (h) show LKFFB tuning failure. KF algorithms set up with $q = 100$ for AKF; $J^{(B)} = 100, \omega_0^{(B)}/2\pi = 0.5$ Hz for LKFFB; with $N_T = 2000, N_P = 100$ steps, $\Delta t = 0.001$ s, $r = 20$; $N.L. = 1\%$.

ment noise. In order to probe this numerically, we perform direct comparisons of filter performance under varying measurement-noise strength for both the AKF and LSF. Since autoregressive coefficients learned in (noisy) environments are re-cast in Kalman form, we test measurement-noise filtering in Kalman frameworks enabled by the design parameter R . In Fig. 5 (a), we plot \tilde{L}_{BR} prediction risk for AKF and LSF as a ratio such that a value greater than unity implies LSF out-

performs AKF. In cases (i)-(iv), we increase the applied noise level to our data-sets $\{y_n\}$ representing simulated measurements on f . For applied measurement noise level $N.L. > 1\%$ in (ii)-(iv), we find that $AKF/LSF < 1$ and AKF outperforms LSF for the conditions studied here, with a general trend towards increasing benefits as noise increases until the noise becomes so large (iv) that the benefits fluctuate as a function of n . Calculations of the ensemble-averaged \tilde{L}_{BR} in Fig. 5 (b) demonstrate that

all ratios reported in (a) correspond to a useful forward prediction horizon.

In machine learning or optimal control settings, the robustness of the learning procedure to small changes in the underlying system is an essential characteristic of the algorithm. In our case, we have already seen that the quality of projection of the true dynamics of f onto the LKFFB basis can have a significant impact on the quality of learning and predictive estimation. We explore this initial finding in more detail.

In Fig. 6, we simulate various learning conditions including (a) perfect learning in LKFFB; (b) imperfect projection relative to the LKFFB basis; (c) imperfect projection combined with finite algorithm resolution; and (d) imperfect learning and undersampling relative to true noise bandwidth. The ordering of figure presentation highlights the degree of impact of the introduced pathologies on LKFFB. By contrast we find reasonable model robustness in AKF/LSF at the expense of performance in the somewhat unrealistic perfect learning case.

We expose the underlying optimisation results for choosing an optimal (σ^*, R^*) for LKFFB in Fig. 6 (e)-(h) and for AKF in Fig. 6 (i)-(l). Individual sample points are highlighted as solid dots while low-loss pairs in this 2D space are highlighted for giving low state-estimation [purple] or prediction [crimson] risk via shaded circles. As the model pathologies indicated above increase, these data demonstrate a divergence between regions of the optimisation space which permit low-loss state estimation and forward prediction for LKFFB. In contrast, overlap of low loss Bayes Risk regions do not change for AKF across Fig. 6 (i)-(l).

Kalman filtering algorithms employed here combine recursive state estimation with the establishment of a dynamical model in the Fourier domain. Therefore, one way to explore algorithmic performance is to look directly at the efficacy of spectral estimation relative to the true (here numerically engineered) hidden dynamics of f . For both the LKFFB and AKF we plot the extracted power spectral density, $S(\omega)$, as a function of angular frequency ω , for different measurement sampling conditions in Fig. 7 against the true spectrum used to define f . These simulated experimental conditions match those introduced in Fig. 4 (b).

In the case of LKFFB, we plot the learned instantaneous amplitudes from a single run [blue markers] and for AKF we extract optimised algorithm parameters as described above [red markers]. Under the assertion that the LSF implements an $\text{AR}(q)$ process, the set of trained parameters, $\{\{\phi_{q' \leq q}\}, \sigma^2\}$ from AKF allows us to derive experimentally measurable quantities, including the power spectral density of the dephasing process: $S(\omega) = \sigma^2 \left(2\pi \left| 1 - \sum_{q'=1}^q \phi_{q'} e^{-i\omega q'} \right|^2 \right)^{-1}$ [35].

The critical feature in these data-sets is the existence of a flat-top spectrum possessing a sharp high frequency cutoff. Both classes of Kalman filtering algorithm successfully identify this structure and locate this high-

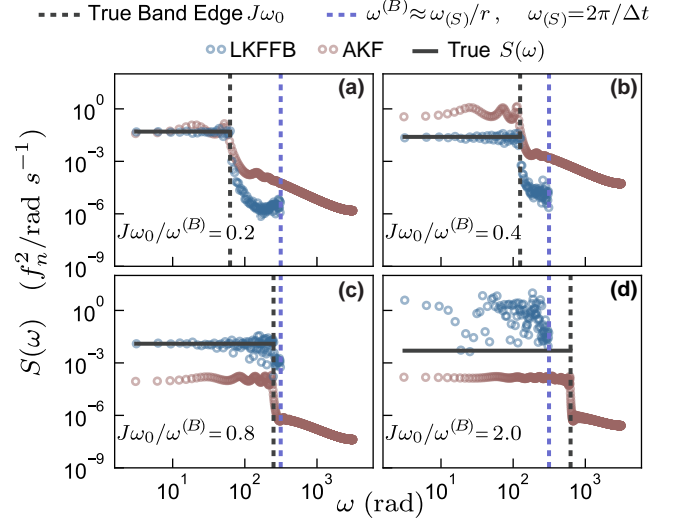


FIG. 7. (a)-(d) Blue (red) open markers plot LKFFB (AKF) spectrum estimates; true spectrum (flat top) of f plotted in black solid line. Dashed black vertical line marks true noise cutoff, $J\omega_0$, and this is varied relative to a measurement sampling rate, $\omega_{(S)}$, and $\omega^{(B)} \equiv \omega_0^{(B)} J^{(B)} \approx \omega_{(S)}/r$ in LKFFB; such that $\omega_0/2\pi = 0.497$ Hz, $J = 20, 40, 80, 200$. For LKFFB, blue open markers are $\propto \|\hat{x}_n^j\|^2$ in a single run with $\omega_0^{(B)}/2\pi = 0.5$ Hz for $j \in J^{(B)} = 100$ oscillators; dashed blue vertical line marks edge of LKFFB basis. For AKF, red markers are $\hat{S}(\omega)$ computed using learned $\{\phi_{q' \leq q}\}$ and optimised σ^* , with order $q = 100$. In all plots, the zeroth Fourier component is omitted on the log scale; and $N_T = 2000$, $N_P = 50$ steps, $\Delta t = 0.001$ s, $r = 20$, with $M = 50$ runs, $K = 75$ trials; $N.L. = 1\%$.

frequency cutoff. In general, however, the LKFFB provides superior spectral estimation relative to the AKF, and enables better estimation of the signal strength in the Fourier domain even in the presence of imperfect projection of f onto the basis used in LKFFB. The only case in which the LKFFB fails is in Fig. 7(d), where the LKFFB basis is ill-specified relative to the true noise bandwidth. The observed behavior is somewhat surprising given the generally superior performance of the AKF in predictive estimation, but does highlight the practical difference between Fourier-domain spectral estimation and time-domain prediction.

C. Performance of the quantised Kalman filter

The discrete nature of projective measurement outcomes in quantum systems poses a potential challenge for Kalman filters in the event that measurement preprocessing as in Fig. 1(b) is not performed. We test filter performance for predictive estimation when only binary measurement outcomes are available via the QKF. To reiterate, QKF estimates and tracks hidden information, f_n , using the Kalman true state x_n . In our construction

the associated probability for a projective qubit measurement outcome, $\propto z_n$ is not inferred or measured directly but given deterministically by Born's rule encoded in the non-linear measurement model, $z_n = h(f_n)$. The measurement action is completed by performing a biased coin flip, where z_n determines the bias of the coin.

For QKF, the normalised ensemble-averaged prediction risk, $\langle (z_n - \hat{z}_n)^2 \rangle_{f,D} / \langle (z_n - \mu_z)^2 \rangle_{f,D}$, is calculated with respect to z as the relevant quantity parameterising qubit-state evolution, instead of the stochastic underlying f . This quantity is labeled as Norm. Risk in Fig. 8 and we test if $\langle (z_n - \hat{z}_n)^2 \rangle_{f,D} / \langle (z_n - \mu_z)^2 \rangle_{f,D} < 1$ for $0 < n < n^*$ can be achieved for numerical experiments considered previously in the linear regime. In particular, we generate true f defined in numerical experiments in Fig. 4(b) (and Fig. 7) for $q = 100$ and varying sample rates.

We isolate the role of the measurement action by first inputting into the QKF a true dynamical model rather than a dynamical model learned as in the standard AKF. To specify true dynamics, we begin with a set of $\{\phi_{q' \leq q}\}$ and exactly derive a new f' . As a result the full set of parameters relevant to the filter, $\{\{\phi_{q' \leq q}\}, \sigma, R\}$, are perfectly defined and known, and the filter simply acts on single shot qubit measurements. These simulations reveal that subject to generic measurement oversampling conditions introduced above the QKF is able to successfully enable predictive estimation. As in the linear case, the absolute forward prediction horizon is arbitrary relative to $\omega_0 J / \omega^{(B)}$ and implicitly, an optimisation over the choice of q for a finite data size, N_T , in our application.

Our simulations reveal that the QKF is considerably more sensitive to measurement noise, model errors, and the degree of undersampling than the linear model as shown in Fig. 8 (b). Here the QKF incorporates a learned dynamical model from AKF in the linear regime and we tune (σ, R) for use in the QKF. In particular, we explore $\sigma \geq \sigma_{AKF}^*$ to incorporate model errors as $\{\phi_{q' \leq q}\}$ were learned in the linear regime. We also incorporate increased measurement noise via $R \geq R_{AKF}^*$ as QKF receives raw data that has not been pre-processed or low-pass filtered. The underlying optimisation problems are well behaved for all cases in Fig. 8(b) [not shown]. As the sampling rate is reduced, the QKF forward prediction horizon collapse rapidly i.e. $\langle (z_n - \hat{z}_n)^2 \rangle_{f,D} / \langle (z_n - \mu_z)^2 \rangle_{f,D} > 1$ prediction risk for all $n > 0$.

D. Failure of GPR in predictive estimation

Under a GPR framework, we test whether predictive performance can be improved by considering the entire measurement record (at once) and projecting this record on an infinite basis of oscillators summarised by a periodic kernel. We investigate several different types of GPR models for $M = 50$ realisations of f in the top panel of Fig. 9. For the results shown, we use a popular choice of a maximum-likelihood optimisation procedure

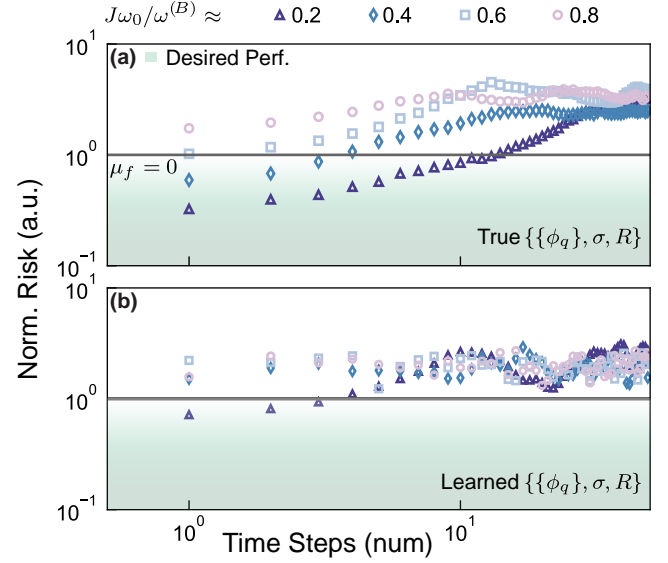


FIG. 8. Norm. Risk against $n > 0$ plotted for QKF in open markers; dark-grey line at $\mu_f \equiv 0$ depicts performance under predicting the noise mean. QKF outperforms predicting the mean if open markers lie in green regions. Marker colour (dark indigo to pink) depicts true noise cutoff varied $J\omega_0/\omega^{(B)} = 0.2, 0.4, 0.6, 0.8$ for f defined identically in Fig. 7 with $\omega_0/2\pi = 0.497$ Hz, $J = 20, 40, 60, 80$; $N.L. = 1\%$. (a) We obtain $\{\phi_{q' \leq q}\}$, $q = 100$ coefficients from AKF/LSF acting on a linear measurement record generated from true f . A new truth, f' , is generated from an $AR(q)$ process using $\{\phi_{q' \leq q}\}$, $q = 100$ as true coefficients and by defining a known, true σ . Quantised measurements from f' are obtained; data is corrupted by measurement noise of a true, known strength R . (b) We use $\{\phi_{q' \leq q}\}$, $q = 100$ coefficients from (a) but we generate quantised measurements from the original, true f . QKF noise design parameters are optimised for ($\sigma_{AKF}^* \leq \sigma_{QKF}$, $R_{AKF}^* \leq R_{QKF}$) with $M = 50$ runs, $K = 75$ trials. For (a)-(b), $N_T = 2000$, $N_P = 50$ steps, $\Delta t = 0.001$ s, $r \gg 2$.

implemented via L-BFGS in GPy [47].

We find that the underlying optimisation procedure for training on our measurement records remains difficult despite having access to an analytical calculation for the cost function. For all results in Fig. 9(a) and (b), we use significant manual tuning prior to deploying the automated procedures in GPy. Hence, we focus on using numerical results under GPR to illuminate structural implications of the choice of kernels in our application, rather than making comparative statements about kernel performance.

The results we have assembled demonstrate that the implementation of GPR with a periodic kernel critically depends on the frequency basis comb spacing, $\omega_0^{(B)}$, or equivalently, a deterministic quantity, κ :

$$\kappa \equiv \frac{2\pi}{\Delta t \omega_0^{(B)}} - N_T \quad (46)$$

The term $2\pi/\Delta t \omega_0^{(B)}$ is the theoretical number of mea-

measurements that, in principle, would be required to *physically* achieve the Fourier resolution set by the kernel hyper-parameter, $\omega_0^{(B)}$, and the fundamentally discrete nature of a sequential Ramsey measurement record, expressed by Δt . Hence, if $\kappa = 0$, the physical Fourier resolution determined by the data set matches the comb spacing in the periodic kernel. For $\kappa > 0$, the comb spacing in the periodic kernel is less than the Fourier spacing defined by the experimental data collection protocol, with total measurements N_T .

In Fig. 9(a), we see that GPR predictive performance for the periodic kernel improves as the Kernel's comb spacing is reduced. For each value of κ we plot \tilde{L}_{BR} against time-steps forward, n^\ddagger , where the \ddagger corresponds to the evaluation of a predictive GPR distribution on arbitrarily chosen test points, $n^\ddagger = -N_T, \dots, -1, 0, 1, \dots, N_P$. Here, the optimiser is constrained to a region in $2\pi/\omega_0^{(B)}$ parameter space that corresponds to the order of magnitude for κ . Grey markers correspond to $\kappa \leq 0$, where the algorithm operates above (or at) the Fourier resolution. In this physically motivated parameter regime, prediction fully fails. It is not until we set $\kappa \sim 10^3$ – a nominally unphysical operating regime where the algorithm's frequency-comb spacing is smaller than the Fourier resolution – that prediction succeeds [red traces]. This latter case is physically difficult to interpret given that in this regime we find the best ensemble-averaged predictive performance only by providing unphysical freedom to the algorithm. We note that the optimised length scale for the periodic kernel remains on of order $\Delta t \sim 10\Delta t$, such that for all red trajectories in panel (a), we are operating in a high $2\pi/\omega_0^{(B)}$, low l limit.

We contextualise the predictive performance of the GPR periodic kernel (PER) [red solid line] in the high- κ , low- l limit by comparing against predictions derived using other standard kernels [dotted lines] in the inset to Fig. 9(a). In such circumstances the predictive performance of the periodic kernel predictive is on par with an application of a Gaussian kernel (RBF) and a scale mixture of zero mean Gaussians with different decay lengths (RQ). A Matern kernel (MAT32) and a quasi-periodic kernel (QPER) yield lower-than-anticipated performance. Further discussion of the choice of kernel appears in Sec. V. For each individual time-trace contributing to the ensemble averages appearing here, we observe that all kernels (PER, RBF, RQ, MAT32, QPER) yield good state estimation and the state estimate at $n^\ddagger = -1$ agrees well with the truth. For GPR with a PER, RBF, and RQ kernels, the state estimate at $n^\ddagger = -1$ smoothly decays to the mean value (zero) for $n^\ddagger \geq 0$ and this effect yields a favourable normalised Bayes prediction risk immediately after $n^\ddagger > 0$ depicted by the solid lines in inset of Fig. 9(a).

In order to illustrate the operating mechanism for the periodic kernel, we dramatically simplify the model used for f in Fig. 9(a) and replace it with a single-frequency sine curve. Fig. 9 (b)-(e) demonstrates the prediction

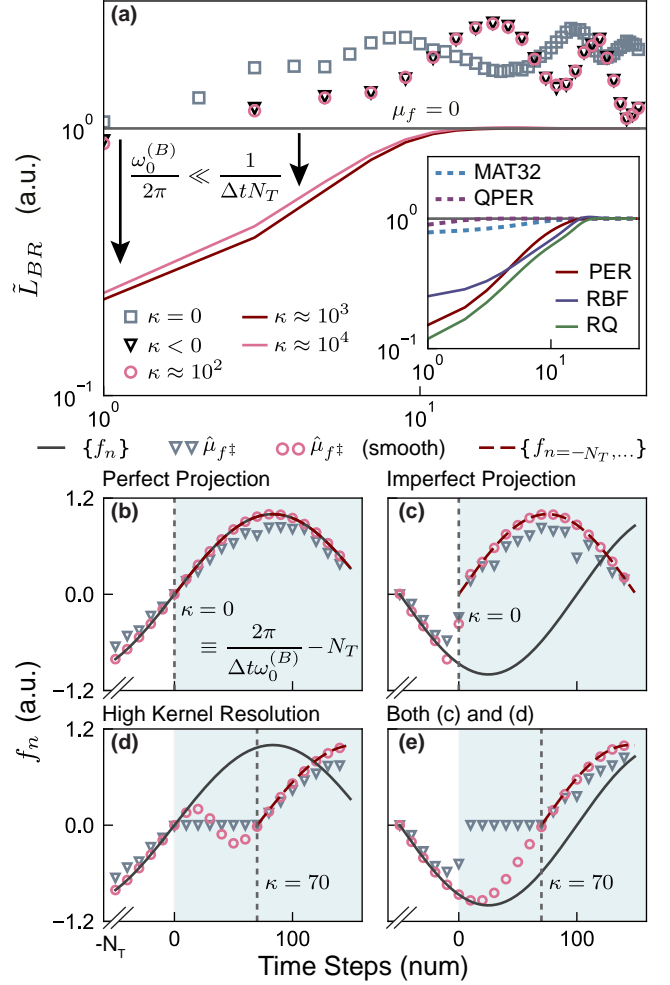


FIG. 9. (a) \tilde{L}_{BR} v. n^\ddagger (in units of number of time steps) are plotted for GPR with a periodic kernel. Dark-grey horizontal line at unity for $\mu_f \equiv 0$ marks \tilde{L}_{BR} under predicting the mean; GPR outperforms predicting the mean if data falls below this line. Grey-black markers correspond to optimisation within physical bounds for $\kappa \leq 0$ (kernel resolution at or above Fourier resolution); crimson markers and lines depict optimisation within unphysical regimes, $\kappa > 0$; with solid lines in high $\kappa \gg 0$ regime. Remaining $\{R, \sigma, l\}$ optimised for non-negative values. Inset (a) \tilde{L}_{BR} v. n^\ddagger of periodic kernel (PER) with $\kappa \approx 10^3$ is plotted against results from naively trained Gaussian kernels (RBF, RQ); a Matern kernel (MAT32) and a quasi-periodic kernel (QPER). (b)-(d) True state f_n v. n [black solid line] and GPR predictions $\hat{\mu}_{f^\ddagger}$ v. n^\ddagger [open markers] plotted for periodic kernel for tracking a sinusoid with frequency, ω_0 ; noisy data record [not shown] ceases at $n = 0$. We fix $\kappa = 0, 70$; triangles plot predictions for manually tuned $\{R, \sigma, l\}$; circles plot predictions for optimised $\{R, \sigma, l\}$. Vertical dashed lines mark $n = \kappa$, where we overlay true f at the beginning of the data record as a red dashed line. (b) Perfection projection is possible $\omega_0/\omega_0^{(B)} \in \mathbb{Z}$ (natural numbers), $\omega_0/2\pi = 3$ Hz. (c) Imperfect projection, with $\omega_0/\omega_0^{(B)} \notin \mathbb{Z}$, $\omega_0/2\pi = 3\frac{1}{3}$ Hz, $\kappa = 0$. (d) Moderately raise $\kappa > 0$, such that $\omega_0/\omega_0^{(B)} \gg 0 \notin \mathbb{Z}$ for original $\omega_0/2\pi = 3$ Hz. (e) Test (c) and (d) for $\kappa > 0$, $\omega_0/\omega_0^{(B)} \notin \mathbb{Z}$, $\omega_0/2\pi = 3\frac{1}{3}$ Hz. For (b)-(e), $N_T = 2000, N_P = 150$ steps, $\Delta t = 0.001$ s; $N.L. = 1\%$.

routine for GPR using a periodic kernel on a simplified version of f , and as before, prediction is always conducted from time-step zero. For this simple example, the periodic kernel learns Fourier information in the measurement record enabling interpolation using test-points $n^\dagger \in [-N_T, 0]$ for all cases (b)-(e) in Fig. 9, and atypical features are seen only for test-points in the prediction region [blue shaded region]. We consider predictions from a manually tuned model [triangles] and an optimised GPR model where remaining free $\{\sigma, R, l\}$ parameters are tuned using GPy [circles].

An examination of different cases for imperfect learning reveal that this discontinuity exhibits deterministic behavior linked to the underlying structure of the algorithm, namely, to the value of κ . In our numerical experiments, we find that in all cases of imperfect learning under GPR with a periodic kernel, a discontinuity in the prediction sequence arises at $n^\dagger = \kappa$. This is marked by the vertical dashed lines in all panels of Fig. 9(b)-(e). However, another feature appears which we identify as being linked to oversampling of the underlying process determining f . In such cases, the algorithm simply predicts zero out to $n^\dagger = \kappa$ before discontinuously predicting future evolution which does not appear similar to the true value of f . By contrast an optimised model gives smoothly varying predictions, which still adhere to the underlying behaviour set by κ for $n^\dagger > 0$.

In Fig. 9(b)-(e), we also plot the value of f as given from $n = -N_T$, the start of the data set, on top of the prediction from $n^\dagger = \kappa$. Here we see that the prediction provided by GPR matches the earliest stages of the underlying data set well. Through various numeric experiments we find that the action of GPR in such parameter regimes (moderately positive $\kappa > 0$) appears to be to simply repeat the learned values of f from $n = -N_T$ beginning at $n^\dagger = \kappa$. Accordingly these predictions rarely describe the underlying forward dynamics of f well.

As we enter the high κ regime, $\kappa \gg 0$, the features in Fig. 9(b)-(e) disappear, and GPR predictions begin to track the (slow moving) ‘truth’ for $n^\dagger \gg 0$. Analogously to inset (a), we see the performance of PER approach that of standard Gaussian kernels in this simplified case.

V. DISCUSSION

The numeric simulations we have performed probe a wide variety of operating conditions in order to explore the algorithmic pathologies of leading forecasting techniques drawn from engineering, econometrics, and machine learning communities when applied to predictive estimation of qubit evolution. A qualitative summary of our observations and key algorithmic differences is given in Table I for ease of reference.

Our central finding is that overall the autoregressive Kalman filter provides an effective path to perform both state estimation and forward prediction for non-Markovian qubit dynamics. Recasting dynamics into

an AKF filter, importantly, provides model robustness against details of the underlying dynamics as well as filtering of noise that allows it to outperform the simpler LSF in [28]. Measurement noise filtering is enabled in the Kalman framework through the optimisation procedure for R and has a regularising (smoothing) effect. Additionally optimisation of the imperfectly learned dynamical model is provided through the tuning of σ . The joint optimisation procedure over (σ, R) ensures that the relative strength of noise parameters is also optimised.

AKF has also been demonstrated to work well with discretised projective measurement models via what we refer to as the QKF. In QKF, we employ single-shot, discretised qubit data while enabling model-robust qubit state tracking and increased measurement noise filtering via the underlying AKF algorithm. However we find that the QKF is vulnerable to the buildup of errors for arbitrary applications and we provide three explanatory remarks from a theoretical perspective. First, the Kalman gains are recursively calculated using a set of linear equations of motion which incorporate the Jacobian H_n of $h(x_n)$ at each n . All non-linear Kalman filters perform well if errors during filtering remain small such that the linearisation assumption holds at all time-steps. Second, measurements are quantised and hence residuals must be $\{-1, 0, 1\}$ rather than continuously represented floating-point numbers. In our case, the Kalman update to x_n at n , mediated by the Kalman gain cannot benefit from a gradual reduction in residuals. A third effect incorporates consequences of both quantised residuals and a non-linear measurement action. In linear Kalman filtering, Kalman gains can be pre-calculated in advance of the acquisition of any measurement data: the recursion of Kalman state-variances P_n , can be decoupled from the recursion of Kalman state-means, x_n [31]. In our application, quantised residuals affect the Kalman update of x_n , and further, they affect the recursion for the Kalman gain via the state dependent Jacobian, H_n .

In this context, we demonstrate numerically that the QKF achieves a desirable forward prediction horizon when the build of errors during filtering is minimised, for example, by specifying Kalman state dynamics and noise strengths perfectly, and/or by severely oversampling relative to the true dynamics of f . At present, we simply interpret our results on the QKF as demonstration that one may in principle track stochastic qubit dynamics using single shot measurements under a Kalman framework. The QKF also has the benefit, as constructed, of reverting to the AKF if suitable pre-processing of data is performed prior to execution of the iterative state-estimation algorithm. In common laboratory settings the measurement protocol may be effectively linearised through simple averaging of multiple single-shot measurements, application of Bayesian estimation protocols, or other pre-processing as identified above. So long as the pre-processing takes place on timescales fast relative to the underlying qubit dynamics, the measurement linearization has no impact other than to change the ef-

Algorithm	Structure	State Est. Perf.	Prediction Perf.	Advantages	Weaknesses
Kalman, AKF	Recursive; autoregressive dynamical model	Good	Best	Robust to measurement noise & variety of operating regimes	Need to train AR model prior to filtering and prediction
Kalman, LKFFB	Recursive; Fourier synthesized dynamical model	Good	Moderate	Robust to measurement noise	Oscillator structure not robust in all operating regimes
Kalman, QKF	Recursive; single qubit data, autoregressive dynamical model	Moderate	Moderate	Direct processing of single shot qubit data	Susceptible to rapid error accumulation via model nonlinearities and binary data
Least Squares, LSF	Batch processing; linear regression	Good	Good	Rapid extraction of autoregressive dynamics from large datasets	Not robust against measurement noise
GPR (PER)	Batch processing; Bayesian data constrained model selection	Good	Poor	Good pattern interpolation during state estimation	Susceptible to producing numeric artifacts in forward prediction

TABLE I. Overview of performance results for all algorithms in this study across all frameworks. Column 2 lists mechanisms for data input (recursive or batch) and the key structural comparisons being made between algorithms. Columns 3-4 qualitatively assess performance during qubit state estimation ($n < 0$) and prediction ($n \geq 0$); we comment on conditions in which algorithms are found to perform strongly or fail in columns 5-6.

fective sample rate of the measurements. Thus it is our view that full implementation of the QKF is not essential if improved optimization routines are not accessible.

It is possible that QKF forward prediction horizons in realistic learning environments can be improved by solving the full $q + 2$ optimisation problem for $\{\{\phi_{q' \leq q}\}, \sigma, R\}$, rather than employing the approach taken in this manuscript. However, this poses its own challenges given the observations we make about the optimisation landscape even for the 2D optimisation problem faced in the AKF. More sophisticated, data-driven model selection schemes are described for both KF and kernel learning machines (such as GPR) in literature (e.g. [48, 49]). Beyond standard local-gradient and simplex optimisers, we consider coordinate ascent [50] and particle swarm optimisation techniques [51] as promising, nascent candidates and their application remains an open research question. One may also consider switching from a high order $\text{AR}(q)$ to an ARMA model with a smaller number of optimisation parameters. Typically, this is accomplished by incorporating either greater a priori information about the underlying dynamic process in the design of the ARMA model and/or using model-less particle-based / unscented filtering techniques to overcome non-linearities in an ARMA representation (e.g. [2]). The latter set of techniques are well adapted for non-linear models but are likely to require a modification to allow for non-Markovian dynamics (e.g. by designing an appropriate transition probability for otherwise Markov re-sampling procedures); in contrast, a typical recursive ARMA formulation for our application may track temporal correlations but be ill-equipped for non-linear, coin-flip measurements. One expects that a straightforward application of such procedures will be complicated.

Our general results on the use of autoregressive models for building Kalman dynamical models stand in contrast to Fourier-domain approaches in LKFFB and GPR us-

ing a periodic kernel; both show significant performance degradation in cases when learning of state dynamics was imperfect. In investigating the loss of performance for LKFFB, we find that the efficacy of this approach depends on a careful choice of a *probe* (i.e. a fixed computational basis) for the dynamics of f capturing the effect of dephasing noise on the qubit. In the imperfect learning regime of Fig. 4 and identically, Fig. 7, LKFFB reconstructs Fourier domain information to a high fidelity across a range of sampling regimes but is outperformed by AKF in the time domain (Fig. 4). Since LKFFB tracks instantaneous amplitude and phase information explicitly for each basis frequency, the loss of LKFFB time-domain predictive performance must accrue from difficulty in tracking instantaneous phase, rather than amplitude, information.

While difficulty of instantaneous phase estimation is likely to disadvantage the time-domain predictive performance of LKFFB, our results show that a Fourier-domain approach yields high fidelity reconstructions of power spectral density describing f . These reconstructions appear robust against imperfect projection on the LKFFB oscillator basis even as oversampling is reduced. This suggests that an application of LKFFB outside of predictive estimation could be tested against standard spectral estimation techniques in future work.

The challenge in adapting GPR for the task of time-domain predictive estimation has proved more striking. In our numerical simulations, under conditions comparable to those tested in the AKF, the values of normalised Bayes prediction risk for all GPR models are at least an order of magnitude greater than the comparable performance of the AKF or LKFFB (refer to panel Fig. 5(b-ii), equivalently, Fig. 6(c)). This difference is somewhat surprising because in the limit that Γ_n is set to the identity in LKFFB and an infinite basis of oscillators in the periodic kernel is truncated at the finite value, $J^{(B)}$, both

LKFFB and the GPR-PER are formally equivalent to classical Kalman filtering for a collection of $J^{(B)}$ independent state-space resonators [42]. In this limit, the true f is described by theoretically identical covariance functions in both KF and GPR frameworks. While we do not operate in this regime, one would expect predictive capabilities of these two algorithms to be comparable.

In contrast to our observations for the various flavors of KF tested here, we observe that GPR predictions with a periodic kernel are useful for filtering/retrodiction but appear to have limited meaning for forward predictions for time-steps $n = n^\dagger > 0$. In our application, predictive performance of GPR with a periodic kernel for $\kappa = 0$ is shown to yield poor predictive performance over the ensemble average (Fig. 9(a)). For the unexpected regime of $\kappa \gg 0$ and relatively small fixed l , predictive performance improves and the periodic kernel performs similarly to RBF and RQ. In this a high κ and a low l regime, the sin term of the periodic kernel is slowly moving ($\sin(x) \approx x$) and hence the argument of the exponential in the periodic kernel approximates a Gaussian, reducing to an RBF kernel. Our numerical investigations show that an optimised RQ kernel consistently chooses parameter regimes where an RQ also converges to an RBF. For the operating regimes pertinent to our application, it appears that the choice of the periodic, RBF, and RQ kernels will produce theoretically equivalent results for forward predictions of the qubit state. In our analysis, these ‘forward predictions’ simply arise from a smoothed decay of state estimates starting from test-point $n^\dagger = -1$ to the noise mean for test-points $n^\dagger > 0$; and are difficult to interpret compared to their Kalman counterparts.

Our numerical characterisation of the periodic kernel for a simple, noiseless f demonstrates that this kernel learns Fourier domain amplitude information in a way that is better suited for pattern fitting than forward prediction. The predictive time domain sequence of state estimates is repetitive at $n = n^\dagger = \kappa$, and can be interpreted as successful qubit-state predictions only when f is perfectly learned (no discontinuities appear). When learning is imperfect, however, GPR with a periodic kernel is able to learn Fourier amplitudes to provide good retrodictive state estimates for $n^\dagger < 0$, but forward predictions for $n^\dagger > 0$ typically fail. Unlike LKFFB, we believe the periodic kernel does not permit actively extracting and updating phase information for each individual basis oscillators at $n^\dagger = \kappa$. Since phase information can be recast as amplitude information for any fixed-frequency oscillator, one would naively expect that forward predictions can be improved by increasing κ moderately, such that the higher order terms in a series expansion of the sin term are non trivial and $\sin(x) \approx x$ cannot apply. However, any positive value of κ means that we are probing dynamics at frequencies lower than appearing in the data-set. As such, a GPR-PER model predicts zero for $n^\dagger \in [0, \kappa]$, $\kappa > 0$, before reviving at κ . The use of a procedure optimising kernel noise parameters $\{\sigma, R\}$ does not change the behavior as $n^\dagger \rightarrow \kappa$, but does smooth

the discontinuities, as illustrated in Fig. 9(f). In letting $\kappa \gg 0$ (extremely large), we lose the uniqueness of the periodic kernel in summarising an infinite basis of oscillators, and standard Gaussian kernels (e.g. RBF, RQ) are likely to apply.

It is possible that the choice of more complex kernels could enhance forward time series predictions via GPR, but they bring additional complications which thus far remain unresolved in relation to the current application. As one example, our ability to use numerical investigations to inform kernel design is further distorted by the need for a robust optimisation procedure, as illustrated by lower-than anticipated predictive performance observed for QPER. Another class of GPR methods, namely, spectral mixture kernels and sparse spectrum approximation using GPR have been explored in [52, 53]. However, these techniques also require efficient optimisation procedures to learn many unknown kernel parameters, whereas the sine-squared exponential in the periodic kernel is parameterised only by two hyper-parameters. Aside from spectral methods, the generalisation of MAT32 to higher $q + 1/2$ models probes only a subset of all possible AR(q) processes, as the restrictions on autoregressive coefficients in Matern kernels are greater than the general case considered under an AKF in this manuscript. A detailed investigation of the application of such methods for forward prediction beyond pattern recognition and with limited computational resources, remains an area of future investigation.

VI. CONCLUSION

In this manuscript, we provided a detailed survey of machine learning and filtering techniques applied to the problem of tracking the state of a qubit undergoing non-Markovian dephasing via a record of projective measurements. We specifically considered the task of performing predictive estimation: learning dynamics of the system from the measurement record and then predicting evolution forward in time. To accommodate stochastic dynamics under arbitrary dephasing, and without an a priori dynamical model, we chose two Bayesian learning protocols - Gaussian Process Regression (GPR) and Kalman Filtering (KF). All Kalman algorithms predicted the qubit state forward in time better than predicting mean qubit behaviour, indicating successful prediction, though an autoregressive approach to building the Kalman dynamical model demonstrated enhanced robustness relative to Fourier-domain approaches. Forward prediction horizons could be arbitrarily increased for all Kalman algorithms by oversampling the underlying dephasing noise. Our investigations included studies of both linear and non-linear measurement routines and validate the utility of the Kalman filtering framework for both. In contrast, under GPR, we found numerical evidence that this approach enables retrodiction but not forward predictions beyond the measurement record.

There are exciting opportunities for machine learning algorithms to increase our understanding of dynamically evolving quantum systems in real time using projective measurements. Quantum systems coupled to classical spatially or temporally varying fields may benefit from classical algorithms to analyse correlation information and enable predictive control of qubits for applications in quantum information, sensing, and the like. Moving beyond a single qubit, we anticipate that measurement records will grow in complexity allowing us to exploit the natural scalability offered by machine learning for mining large datasets. In realistic laboratory environments, the success of algorithmic approaches will be contingent on robust and computationally efficient algorithmic optimisation procedures as well as the extensions beyond Markovian dynamics studied here. The pursuit of these opportunities is the subject of ongoing research.

VII. ACKNOWLEDGMENTS

The LSF filter is written by V. Frey and S. Mavadia [28]. The GPR framework is implemented and optimised using standard protocols in GPy [47]. Authors thank C. Granade, K. Das, V. Frey, S. Mavadia, H. Ball, C. Ferrie and T. Scholten for useful comments. This work partially supported by the ARC Centre of Excellence for Engineered Quantum Systems CE110001013, the US Army Research Office under Contract W911NF-12-R-0012, and a private grant from H. & A. Harley.

Appendix A: Physical Setting

In this Appendix, we derive Eq. (1). We consider a qubit under environmental dephasing. For any two level system, a quantum mechanical description of physical quantities of interest can be provided in terms of the Pauli spin operators $\{\hat{\sigma}_x, \hat{\sigma}_y, \hat{\sigma}_z\}$. If $\hbar\omega_A$ corresponds to an energy difference separating these two qubit states, then the Hamiltonian for a single qubit in free evolution can be written in the Pauli representation. We consider a qubit states in the $\hat{\sigma}_z$ basis, $|0\rangle$ or $|1\rangle$ with energies E_0, E_1 in our notation, corresponding to a 0 or 1 outcome upon measurement. This yields a Hamiltonian for a single qubit as:

$$\hat{\sigma}_z \equiv |1\rangle\langle 1| - |0\rangle\langle 0| \quad (\text{A1})$$

$$\hat{\mathcal{I}} \equiv |0\rangle\langle 0| + |1\rangle\langle 1| \quad (\text{A2})$$

$$E_{0,1} \equiv \mp \frac{1}{2} \hbar\omega_A \quad (\text{A3})$$

$$\hat{\mathcal{H}}_0 = \frac{1}{2} (E_0 |0\rangle\langle 0| + E_1 |1\rangle\langle 1|) \quad (\text{A4})$$

$$+ \frac{1}{2} [(E_1 - E_0) \hat{\sigma}_z + E_0 |1\rangle\langle 1| + E_1 |1\rangle\langle 1|] \quad (\text{A5})$$

$$= \frac{1}{2} \hbar\omega_A \hat{\sigma}_z \quad (\text{A6})$$

In this representation, the effect of dephasing noise on a free qubit system is that any initially prepared qubit superposition of $|0\rangle$ and $|1\rangle$ states will decohere over time in the presence of dephasing noise. This physical effect is modelled as a stochastically fluctuating process $\delta\omega(t)$ that couples with the $\hat{\sigma}_z$ operator. The noise Hamiltonian is described as:

$$\hat{\mathcal{H}}_N(t) \equiv \frac{\hbar}{2} \delta\omega(t) \hat{\sigma}_z \quad (\text{A7})$$

In the formula above, $\delta\omega(t)$ is a classical, stochastically fluctuating parameter that models environmental dephasing and $\hbar/2$ appears as a convenient scaling factor. The total Hamiltonian for a single qubit under dephasing is:

$$\hat{\mathcal{H}}(t) \equiv \hat{\mathcal{H}}_0 + \hat{\mathcal{H}}_N(t) \quad (\text{A8})$$

Since $\hat{\mathcal{H}}_N(t)$ commutes with $\hat{\mathcal{H}}_0$, we can transform away $\hat{\mathcal{H}}_0$ by moving to a rotating frame with respect to H_0 . Let $|\psi(t)\rangle$ be a state in the lab frame, let \hat{U} define a transformation to a rotating frame, and let $|\tilde{\psi}(t)\rangle$ be the state in the rotating frame. The notation, $\tilde{\cdot}$, indicates operators and states in the transformed frame. In this simple case, the transformed Hamiltonian governing the evolution of $|\tilde{\psi}(t)\rangle$ will just be $\hat{\mathcal{H}}_N(t)$:

$$\hat{U} \equiv e^{-i\hat{\mathcal{H}}_0 t/\hbar} \quad (\text{A9})$$

$$|\tilde{\psi}(t)\rangle \equiv \hat{U}^\dagger |\psi(t)\rangle \quad (\text{A10})$$

$$i\hbar \frac{d}{dt} |\tilde{\psi}(t)\rangle \equiv i\hbar \frac{d}{dt} \hat{U}^\dagger |\psi(t)\rangle \quad (\text{A11})$$

$$= -\hat{\mathcal{H}}_0 \hat{U}^\dagger |\psi(t)\rangle + i\hbar \hat{U}^\dagger \frac{d}{dt} |\psi(t)\rangle \quad (\text{A12})$$

$$= (\hat{U}^\dagger \mathcal{H}(t) \hat{U} - \hat{\mathcal{H}}_0) |\tilde{\psi}(t)\rangle \quad (\text{A13})$$

$$\implies \hat{\mathcal{H}} \equiv \hat{U}^\dagger \mathcal{H}(t) \hat{U} - \hat{\mathcal{H}}_0 \quad (\text{A14})$$

$$= \hat{U}^\dagger \hat{\mathcal{H}}_0 \hat{U} + \hat{U}^\dagger \hat{\mathcal{H}}_N(t) \hat{U} - \hat{\mathcal{H}}_0 \quad (\text{A15})$$

$$= \hat{\mathcal{H}}_N(t), \quad [\hat{U}, \hat{\mathcal{H}}_0] = [\hat{U}, \hat{\mathcal{H}}_N(t)] = 0 \quad (\text{A16})$$

$$(\text{A17})$$

In the semiclassical approximation, $\hat{\mathcal{H}}_N(t)$ commutes with itself at different t , and hence we can write a unitary time evolution operator in the rotating frame as:

$$\hat{\tilde{U}}(t, t + \tau) \equiv e^{-\frac{i}{\hbar} \int_t^{t+\tau} \hat{\mathcal{H}}_N(t') dt'} = e^{-\frac{i}{2} f(t, t+\tau) \hat{\sigma}_z} \quad (\text{A18})$$

$$f(t, t + \tau) \equiv \int_t^{t+\tau} \delta\omega(t') dt' \quad (\text{A19})$$

In the rotating frame, we prepare an initial state that is a superposition of $|0\rangle$ and $|1\rangle$ states. This state evolves under $\hat{\mathcal{H}}_N(t)$ during a Ramsey experiment for duration τ . Subsequently, the qubit state is rotated before a projective measurement is performed with respect to the $\hat{\sigma}_z$ axis i.e. the measurement action resets the qubit.

Without loss of generality, define the initial state as $|\tilde{\psi}(0)\rangle \equiv \frac{1}{\sqrt{2}}|0\rangle + \frac{1}{\sqrt{2}}|1\rangle$ in the rotating frame. Then, the probability of measuring the same state after time τ in a single shot measurement, d_n as:

$$Pr(d_n = 1|f(0, \tau), \tau) = |\langle \tilde{\psi}(0)|\hat{U}(0, \tau)|\tilde{\psi}(0)\rangle|^2 \quad (\text{A20})$$

$$Pr(d_n = 0|f(0, \tau), \tau) \equiv 1 - Pr(d_n = 1|f(0, \tau), \tau) \quad (\text{A21})$$

The second $\pi/2$ control pulse rotates the state vector such that a measurement in $\hat{\sigma}_z$ basis is possible, and the probabilities correspond to observing the qubit in the $|1\rangle$ state. Hence, Eq. (A20) defines the likelihood for single shot qubit measurement. Further, Eq. (A20) defines the non linear measurement action on phase noise jitter, $f(0, \tau)$. We impose a condition that $f(0, \tau)/2 \leq \pi$ such that accumulated phase over τ can be inferred from a projective measurement on the $\hat{\sigma}_z$ axis.

1. Experimentally Controlled Discretisation of Dephasing Noise

In this section, we consider a sequence of Ramsey measurements. At time t , the Eq. (A20) describes the qubit measurement likelihood at one instant under dephasing noise. We assume that the dephasing noise is slowly drifting with respect to a fast measurement action on timescales of order τ . In this regime, Eq. (A19) discretises the continuous time process $\delta\omega(t)$, at time t , for a number of $n = 0, 1, \dots, N$ equally spaced measurements with $t = n\Delta t$. Performing the integral for $\tau \ll \Delta t$ and slowly drifting noise such that we substitute the following terms in Eq. (A19):

$$\delta\bar{\omega}_n \equiv \delta\omega(t')|_{t'=n\Delta t} \quad (\text{A22})$$

$$f_n \equiv f(n\Delta t, n\Delta t + \tau) \quad (\text{A23})$$

$$= \frac{\hbar}{2} \int_{n\Delta t}^{n\Delta t + \tau} \delta\bar{\omega}_n dt' = \frac{\hbar}{2} \hat{\sigma}_z \delta\bar{\omega}_n \tau \quad (\text{A24})$$

In this notation, $\delta\bar{\omega}_n$ is a random variable realised at time, $t = n\Delta t$, and it remains constant over short duration of the measurement action, τ . We use the shorthand $f_n \equiv f(n\Delta t, n\Delta t + \tau)$ to label a sequence of stochastic, temporally correlated qubit phases $f \equiv \{f_n\}$.

Since the qubit is reset by each projective measurement at n , the unitary operator governing qubit evolution is also reset such that $\{\hat{U}_n \equiv \hat{U}(n\Delta t, n\Delta t + \tau)\}$ are a collection of N unitary operators describing qubit evolution for each new Ramsey experiment. They are not to be interpreted, for example, as describing qubit free evolution without re-initialising the system. Hence, for each stochastic qubit phase f_n , the true probability for observing the $|1\rangle$ in a single shot is given by substituting f_n for $f(0, 1)$ in Eq. (A20).

$$Pr(d_n = d|f_n, \tau, n\Delta t) = \begin{cases} \cos(\frac{f_n}{2})^2 & \text{for } d = 1 \\ \sin(\frac{f_n}{2})^2 & \text{for } d = 0 \end{cases} \quad (\text{A25})$$

The last line follows from the fact that total probability of the qubit occupying either state must add to unity. This yields Eq. (1) in the main text.

2. True Dephasing Noise Engineering

In the absence of an a priori model for describing qubit dynamics under dephasing noise, we impose the following properties on a sequence of stochastic phases, $f \equiv \{f_n\}$ such that we can design meaningful predictors of qubit state dynamics. We assert that a stochastic process, f_n , indexed by a set of values, $n = 0, 1, \dots, N$ satisfies:

$$\mathbb{E}[f_n] = \mu_f \quad \forall n \quad (\text{A26})$$

$$\mathbb{E}[f_n^2] < \infty \quad \forall n \quad (\text{A27})$$

$$\mathbb{E}[(f_{n_1} - \mu)(f_{n_2} - \mu)] = R(\nu), \quad \nu = |n_1 - n_2|, \quad \forall n_1, n_2 \in N \quad (\text{A28})$$

$$R(\nu) \neq \sigma^2 \delta(\nu) \quad (\text{A29})$$

Covariance stationarity of f is established by satisfying Eqs. (A26) to (A28), namely that the mean is independent of n , the second moments are finite, and the covariance of any two stochastic phases at arbitrary time-steps, n_1, n_2 , do not depend on time steps but only on the separation distance, ν . The $\delta(\nu)$ in the last condition, Eq. (A29), is the Dirac-delta function and establishes that f is not delta-correlated (white). This condition captures the slowly drifting assumption for environmental dephasing noise.

We also require that correlations in f eventually die off as $\nu \rightarrow \infty$ otherwise any sample statistics inferred from noise-corrupted measurements are not theoretically guaranteed to converge to the true moments. Let M be the number of runs for an experiment with M different realisations of the random process f , μ_f be the true mean, $\hat{\mu}_f$ its estimate, \mathcal{D}_M denote the dataset of M experiments, and $R(\nu)$ define the correlation function for the true process, f . Then mean square ergodicity states that estimators approach true moments only if the correlations die off over long temporal separations:

$$\begin{aligned} \lim_{M \rightarrow \infty} \frac{1}{M} \sum_{\nu=0}^{M-1} R(\nu) = 0 &\iff \lim_{M \rightarrow \infty} \mathbb{E}[(\hat{\mu}_f - \mu_f)^2]_{\mathcal{D}_M} = 0 \\ &\text{for } \nu = |n_{m_1} - n_{m_2}|, \quad \forall m_1, m_2 \in M, n_{m_1}, n_{m_2} \in N \\ &\text{with } \hat{\mu}_f = \frac{1}{M} \sum_{m=0}^M f_{n_m} \end{aligned} \quad (\text{A30})$$

The statement above means that a true $R(\nu)$ associated with f is bandlimited for sufficiently large (but unknown) M . If correlations never ‘die out’, then any designed predictors for one realisation of dephasing noise will fail for a different realisation of the same true dephasing. For the purposes of experimental noise engineering, we satisfy the assumptions above by engineering discretised process, f , as:

$$f_n = \alpha \omega_0 \sum_{j=1}^J j F(j) \cos(\omega_j n \Delta t + \psi_j) \quad (\text{A31})$$

$$F(j) = j^{\frac{\eta}{2}-1} \quad (\text{A32})$$

As described in [46], α is an arbitrary scaling factor, ω_0 is the fundamental spacing between true adjacent discrete frequencies, such that $\omega_j = 2\pi f_0 j = \omega_0 j, j = 1, 2, \dots, J$. For each frequency component, there exists a uniformly distributed random phase, $\psi_j \in [0, \pi]$. The free parameter η allows one to specify an arbitrary shape of the true power spectral density of f . In particular, the free parameters $\alpha, J, \omega_0, \eta$ are true dephasing noise parameters which any prediction algorithm cannot know beforehand.

It is straightforward to show that f is covariance stationary. To show mean square ergodicity of f , one requires phases are randomly uniformly distributed over one cycle for each harmonic component of f [54]. Subsequently, one shows that an ensemble average and a long time average of multi-component engineered f are equal. For the evaluation of the long time average, we use product-to-sum formulae and observe that the case $j \neq j'$ has a zero contribution as any finite contribution from cosine terms over a symmetric integral are reduced to zero as $N \rightarrow \infty$. For $j = j'$, only a single cosine term survives. The surviving term depends on ν and N cancels to yield a finite, non-zero contribution that matches the ensemble average.

We briefly comment that f is Gaussian by the central limit theorem in the regimes considered in this manuscript. The probability density function of a sum of random variables is a convolution of the individual probability density functions. The central limit theorem grants that each element of f_n at n appears Gaussian distributed for large J , irrespective of the underlying properties of the constituent terms or the distribution of the phases ψ . Numerical analysis shows that $J > 15$ results in each f_n appearing approximately Gaussian distributed.

There is an important difference between f_n - defined here in Appendix A and - and f_n in Appendices B and C. In subsequent Appendices B and C, the term f_n defines the ‘true model’ for an algorithmic representation of an arbitrary covariance stationary process - either by invoking Wold’s decomposition theorem (AKF, QKF) or the spectral representation theorem (LKFFB, GPR with Periodic Kernel). This means that f_n in subsequent Appendices only approximates the true covariance stationary stochastic qubit phases, $\{f_n\}$ of the Appendix A in the limit where total size of available sample data increases to infinity. Our notation, f_n , fails to distinguish these two different interpretations as such a difference does not arise in typical applications - in our case, we have no a priori true model of describing stochastic qubit phases, and must rely on mean square approximations. Henceforth, we retain f_n to be the true model for an algorithm with an understanding that this refers to an approximate representation of an arbitrary, covariance stationary sequence of stochastic qubit phases. We reserve the use of the \hat{f}_n for the state estimates and predictions that an algorithm makes having considered a single noisy measurement record.

Appendix B: Autoregressive Representation of f in AKF (and QKF)

Our objective in this Appendix is to justify the representation of f_n assumed by the AKF. In particular, we justify any f_n drawn from any arbitrary power spectral density satisfying the properties in Appendix A 2 can be approximated by a high order autoregressive process.

Such results are well known, if dispersed among standard engineering and econometrics textbooks [4, 11, 33–35, 55]. We struggled to find standard references that explicitly link high q AR models in approximating arbitrary covariance stationary time series of arbitrary power spectral densities, though some general comments are made in [55]. In the discussion below, we summarise relevant background, and link a high q AR process to a theorem that guarantees arbitrary representation of zero mean covariance stationary processes, and provide explicit references for proofs out of scope of introductory remarks in this Appendix. In order to achieve this, we will consider autoregressive (AR) processes of order q , (AR(q)), and moving average processes of order, p (MA(p)). A model incorporating both types of processes is known as an ARMA(q, p) model in our notation.

First, we define the lag operator, \mathcal{L} . This operator defines a map between time series sequences and enables a compact description of ARMA processes. For an infinite time series $\{f_n\}_{n=-\infty}^{\infty}$ and a constant scalar, c , the lag operator is defined by the following properties:

$$\mathcal{L}f_n = f_{n-1} \quad (\text{B1})$$

$$\mathcal{L}^q f_n = f_{n-q} \quad (\text{B2})$$

$$\mathcal{L}(cf_n) = c\mathcal{L}f_n = cf_{n-1} \quad (\text{B3})$$

$$\mathcal{L}f_n = c, \quad \forall n, \implies \mathcal{L}^q f_n = c \quad (\text{B4})$$

Next, we define a Gaussian white noise sequence, ξ , under the strong condition than what is stated simply in Eq. (B6), that ξ_{n_1}, ξ_{n_2} are independent $\forall n_1, n_2$:

$$\mathbb{E}[\xi] \equiv 0 \quad (\text{B5})$$

$$\mathbb{E}[\xi_{n_1}\xi_{n_2}] \equiv \sigma^2\delta(n_1 - n_2) \quad (\text{B6})$$

With these definitions, we can define an autoregressive process and a moving average process of unity order. Eq. (B7) defines an AR($q = 1$) process and dynamics of f_n are given as lagged values of the variable f . The second definition in Eq. (B8) depicts a MA($p = 1$) process where dynamics are given by lagged values of Gaussian white noise ξ .

$$(1 - \phi_1\mathcal{L})f_n = c + \xi_n \quad (\text{B7})$$

$$f_n = c' + (\Psi_1\mathcal{L} + 1)\xi_n \quad (\text{B8})$$

Here, Ψ_1, ϕ_1 are known scalars defining dynamics of f_n ; w_n is a white noise Gaussian process, and c, c' are fixed scalars. It is well known that an MA(∞) representation is equivalently an AR(1) process, and the reverse relationship also applies. For example, we can re-write Eq. (B7) as:

$$f_n = c + \xi_n + \phi_1 f_{n-1} \quad (\text{B9})$$

$$= w_n + \phi_1 f_{n-1} \quad (\text{B10})$$

$$= w_n + \phi_1(w_{n-1} + \phi_1 f_{n-2}) \quad (\text{B11})$$

$$\vdots \quad (\text{B12})$$

$$= \phi_1^{n+1}F_0 + \phi_1^n w_0 + \phi_1^{n-1}w_1 + \dots w_n \quad (\text{B13})$$

$$= \phi_1^{n+1}F_0 + \phi_1^n(c + \xi_0) + \dots + (c + \xi_n) \quad (\text{B14})$$

$$= \phi_1^{n+1}F_0 + c(\phi_1^n + \phi_1^{n-1} + \dots + 1) + \sum_{k=0}^n \phi_1^k \xi_{n-k} \quad (\text{B15})$$

$$w_n \equiv c + \xi_n \quad (\text{B16})$$

$$F_0 \equiv f_{n=-1} \quad (\text{B17})$$

In the last line (and for all subsequent analysis in this Appendix), k should only be interpreted as a index variable for compactly re-writing terms in an equation as summations. We restrict $|\phi_1| < 1$ such that f is covariance stationary [34]. Under these conditions, we take the limit of f capturing an infinite past, namely, as $n \rightarrow \infty$. The initial state

F_0 is eventually forgotten, $\phi_1^{n+1}F_0 \approx 0$ if n is large and $|\phi_1| < 1$. Similarly, the terms $c(\phi_1^n + \phi_1^{n-1} + \dots + 1)$ can be summarised as a geometric series in ϕ_1 . The remaining terms satisfy the definition of an MA(∞) process:

$$f_n = c \frac{1}{1 - |\phi_1|} + \sum_{k=0}^{\infty} \phi_1^k \xi_{n-k}, \quad |\phi_1| < 1 \quad (\text{B18})$$

It is straightforward to show that the reverse is true, namely, an MR(1) is equivalent to an AR(∞) representation [34].

The consideration of an MA(∞) process leads us directly to Wold's decomposition for arbitrary covariance stationary processes, namely, that any covariance stationary f can be represented as:

$$f_n \equiv c' + \sum_{k=0}^{\infty} \Psi_k \mathcal{L}^k \xi_n \quad (\text{B19})$$

$$c' \equiv \mathbb{E}[f_n | f_{n-1}, f_{n-2}, \dots] \quad (\text{B20})$$

$$\Psi_0 \equiv 1 \quad (\text{B21})$$

$$\sum_{k=0}^{\infty} \Psi_k^2 < \infty \quad (\text{B22})$$

Eq. (B19) defines an MA(∞) process derived previously as an AR(1) process. This process is ergodic for Gaussian ξ . However, such a representation of f requires fitting data to an infinite number of parameters $\{\Psi_1, \Psi_2, \dots\}$ and approximations must be made.

We approximate an arbitrary covariance stationary f using finite but high order AR(q) processes. Below we show that any finite order AR(q) process has an MA(∞) representation satisfying Wold's theorem.

We define an arbitrary AR(q) process as:

$$\xi_n \equiv (1 - \phi_1 \mathcal{L} - \phi_2 \mathcal{L}^2 - \dots - \phi_q \mathcal{L}^q)(f_n - c) \quad (\text{B23})$$

In particular, we define $\lambda_i, i = 1, \dots, q$ as eigenvalues of the dynamical model, Φ :

$$\Phi \equiv \begin{bmatrix} \phi_1 & \phi_2 & \phi_3 & \dots & \phi_{q-1} & \phi_q \\ 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 & 0 \end{bmatrix} \quad (\text{B24})$$

$$\lambda \equiv [\lambda_1 \dots \lambda_q] \quad \text{s.t.} |\Phi - \lambda \mathcal{I}_q| = 0 \quad (\text{B25})$$

$$(\text{B26})$$

We use the following result from [34] without proof that the above implies:

$$1 - \phi_1 \mathcal{L} - \phi_2 \mathcal{L}^2 - \dots - \phi_q \mathcal{L}^q \quad (\text{B27})$$

$$\equiv (1 - \lambda_1 \mathcal{L}) \dots (1 - \lambda_q \mathcal{L}) \quad (\text{B28})$$

This yields:

$$\xi_n = (1 - \lambda_1 \mathcal{L}) \dots (1 - \lambda_q \mathcal{L})(f_n - c) \quad (\text{B29})$$

For us to invert this problem and recover an MA process, we need to show that the inverse for each $(1 - \lambda_{q'} \mathcal{L})$ term exists for $q' = 1, \dots, q$. To do this, we start by defining the operator $\Lambda_q(\mathcal{L})$:

$$\Lambda_q(\mathcal{L}) \equiv \lim_{k \rightarrow \infty} (1 + \lambda_q \mathcal{L} + \dots + \lambda_q^k \mathcal{L}^k) \quad (\text{B30})$$

We consider an arbitrary q' -th eigenvalue term in process and we multiply $\Lambda_{q'}(\mathcal{L})$:

$$\Lambda_{q'}(\mathcal{L})\xi_n = \Lambda_{q'}(\mathcal{L})(1 - \lambda_0\mathcal{L}) \dots (1 - \lambda_{q'}\mathcal{L}) \dots (f_n - c) \quad (\text{B31})$$

$$= \lim_{k \rightarrow \infty} (1 + \lambda_{q'}\mathcal{L} + \dots + \lambda_{q'}^k\mathcal{L}^k)(1 - \lambda_{q'}\mathcal{L})(1 - \lambda_0\mathcal{L}) \dots (1 - \lambda_{q'-1}\mathcal{L})(1 - \lambda_{q'+1}\mathcal{L}) \dots (1 - \lambda_q\mathcal{L})(f_n - c) \quad (\text{B32})$$

$$= \lim_{k \rightarrow \infty} (1 + \lambda_{q'}\mathcal{L} + \dots + \lambda_{q'}^k\mathcal{L}^k)(1 - \lambda_0\mathcal{L}) \dots (1 - \lambda_{q'-1}\mathcal{L})(1 - \lambda_{q'+1}\mathcal{L}) \dots (1 - \lambda_q\mathcal{L})(f_n - c) \quad (\text{B33})$$

$$= \lim_{k \rightarrow \infty} (\lambda_{q'}\mathcal{L} + \dots + \lambda_{q'}^{k+1}\mathcal{L}^{k+1})(1 - \lambda_0\mathcal{L}) \dots (1 - \lambda_{q'-1}\mathcal{L})(1 - \lambda_{q'+1}\mathcal{L}) \dots (1 - \lambda_q\mathcal{L})(f_n - c) \quad (\text{B34})$$

$$= \lim_{k \rightarrow \infty} (1 + \lambda_{q'}^{k+1}\mathcal{L}^{k+1})(1 - \lambda_0\mathcal{L}) \dots (1 - \lambda_{q'-1}\mathcal{L})(1 - \lambda_{q'+1}\mathcal{L}) \dots (1 - \lambda_q\mathcal{L})(f_n - c) \quad (\text{B35})$$

Each of the residual terms, $\lambda_{q'}^{k+1}\mathcal{L}^{k+1} \rightarrow 0$ if $|\lambda_{q'}| < 1$ for large k , and this case $\Lambda_{q'}(\mathcal{L})$ defines the inverse $(1 - \lambda_{q'}\mathcal{L})^{-1}$. This procedure is repeated for all q eigenvalues to invert Eq. (B29) and subsequently perform a partial fraction expansion as follows:

$$f_n - c = \frac{1}{(1 - \lambda_1\mathcal{L}) \dots (1 - \lambda_q\mathcal{L})} \xi_n \quad (\text{B36})$$

$$= \sum_{q'=1}^q \frac{a_{q'}}{1 - \lambda_{q'}\mathcal{L}} \xi_n \quad (\text{B37})$$

$$a_{q'} \equiv \frac{\lambda_{q'}^{q-1}}{\prod_{q''=1, q'' \neq q'}^q (\lambda_{q'} - \lambda_{q''})} \quad (\text{B38})$$

The coefficients are $a_{q'}$ as obtained via the partial fraction expansion method during which \mathcal{L} is treated as an ordinary polynomial. At present, we have a represent f via a finite q weighted average of values of ξ . However, in substituting the definition of $\Lambda_{q'} \equiv (1 - \lambda_{q'}\mathcal{L})^{-1}$ from Eq. (B30) in Eq. (B37) and regrouping terms in powers of \mathcal{L} , we recover the form of an MA representation (setting $c \equiv \bar{f}_n = 0$, $\forall n$ for simplicity):

$$f_n = \left[\sum_{q'=1}^q a_{q'}\mathcal{L}^0 + \lim_{k \rightarrow \infty} \sum_{k'=1}^k \left(\sum_{q'=1}^q a_{q'}\lambda_{q'}^{k'} \right) \mathcal{L}^{k'} \right] \xi_n \quad (\text{B39})$$

$$= \Psi_0 + \sum_{k=1}^{\infty} \Psi_k \mathcal{L}^k \xi_n \quad (\text{B40})$$

$$\Psi_0 \equiv \sum_{q'=1}^q a_{q'}\mathcal{L}^0 \quad (\text{B41})$$

$$\Psi_k \equiv \sum_{q'=1}^q a_{q'}\lambda_{q'}^{k'} \quad (\text{B42})$$

By examining the properties of Φ raised to arbitrary powers, it can be shown that $\sum_{q'=1}^q a_{q'} \equiv 1$ and Ψ_k is the first element of Φ raised to the k -th power [34], yielding absolute summability of Ψ_k if $|\phi_{q' < q}| < 1$. This ensures that Wold's theorem is fully satisfied and an $\text{AR}(p)$ process has an $\text{MA}(\infty)$ representation. In moving to an arbitrarily high q , we enable the approximation of any covariance stationary f .

The proofs that high q AR approximations for covariance stationary f improve with q for example, in [37]. The key correspondence is that the number of finite lag terms q in an $\text{AR}(q)$ model contribute to the first q values of the covariance function. This approximation improves with q even if f is not a true AR process [37, 55]. Asymptotically efficient coefficient estimates for any $\text{MA}(\infty)$ representation of f are obtained by letting the order of a purely $\text{AR}(q)$ process tend to infinity and increasing total data size, N [37].

When data is fixed at N , we expect a high q model to gradually saturate in predictive estimation performance. One can arbitrarily increase performance by increasing both q, N [37]. In our application with finite data N , we increase q to settle on a high order AR model while training LSF to track arbitrary covariance stationary power spectral densities [35].

A high q AR model is often the first step for developing models with smaller number of parameters, for example, considering a mixture of finite order $\text{AR}(q)$ and $\text{MA}(p)$ models and estimating $p + q$ number of coefficients using a range of standard protocols [35, 55]. The design of potential ARMA models for our application requires further investigation beyond the scope of this manuscript.

Appendix C: Spectral Representation of f in GPR (Periodic Kernel) and LKFFB

The well-known spectral representation theorem guarantees that any covariance stationary random process (real or complex) can be represented in a generalised harmonic basis. We defer a detailed treatment of spectral analysis of covariance stationary processes in standard textbooks, for example, [34, 38] and present background and key results to provide insights into the choice of LKFFB and GPR (periodic kernel).

The spectral representation theorem states that any covariance stationary random process has a representation given by f_n , and correspondingly, a probability distribution, $F(\omega)$ over $[-\pi, \pi]$ in the dual domain such that:

$$f_n = \mu_f + \int_0^\pi [a(\omega) \cos(\omega n) + b(\omega) \sin(\omega n)] d\omega \quad (C1)$$

$$R(\nu) = \int_{-\pi}^\pi e^{-i\omega\nu} dF(\omega) \quad (C2)$$

Here, μ_f is the true mean of the process f . The processes $a(\omega)$ and $b(\omega)$ are zero mean and serially and mutually uncorrelated, namely, $\int_{\omega_1}^{\omega_2} a(\omega) d\omega$ is uncorrelated with $\int_{\omega_3}^{\omega_4} a(\omega) d\omega$ and $\int_{\omega_j}^{\omega_{j'}} b(\omega) d\omega$ for any $\omega_1 < \omega_2 < \omega_3 < \omega_4$ and any choice of j, j' within the half cycle $[0, \pi]$.

The distribution $F(\omega)$ exists as a limiting case of considering cumulative probability density functions for f_n at each n and letting $n \rightarrow \infty$ such that a sequence of these density functions approach $F(\omega)$ [38]. If $F(\omega)$ is differentiable with respect to ω , then we see the power spectral density $S(\omega)$ and $R(\nu)$ are Fourier duals [38]:

$$R(\nu) = \int_{-\pi}^\pi e^{-i\omega\nu} S(\omega) d\omega \quad (C3)$$

$$S(\omega) \equiv \frac{dF(\omega)}{d\omega} \quad (C4)$$

The duality of the covariance function and the spectral density is formally expressed in literature by the Wiener Khinchin theorem.

We consider the finite sample analogue of the spectral representation theorem considered above by following [34]. To proceed, we define mean square convergence as a distance metric for determining when a sequence of random variables $\{\hat{f}_n\}$ converges to a random variable, f_n in the mean square limit if:

$$\mathbb{E}[\hat{f}_n^2] < \infty \quad \forall n \quad (C5)$$

$$\lim_{n \rightarrow \infty} \mathbb{E}[\hat{f}_n - f_n] = \lim_{n \rightarrow \infty} \|\hat{f}_n - f_n\| = 0 \quad (C6)$$

The statement $\|\hat{f}_n - f_n\| = 0$ measures the closeness between random variables \hat{f}_n and f_n even though the mean square limit is defined for terms of a sequence of random variables, $\{\hat{f}_n\}$, where convergence improves with $n \rightarrow \infty$. In context of this study, we define \hat{f}_n as a linear predictor of f_n belonging to a covariance stationary f . Hence, each \hat{f}_n for large n is a linear combination of the set of random variables belonging all past noisy observations (and in Kalman Filtering, all past state predictions). Mean square convergence of $\|\hat{f}_n - f_n\| = 0$ in our context is a statement of the quality of a predictor, \hat{f}_n , in predicting f_n as the total measurement data grows.

Next, we account for finite data and define the finite sample analogue for the spectral representation theorem. We suppose there exists a set of arbitrary, fixed frequencies $\{\omega_j\}$ for $j = 1, \dots, J$. We let n denote finite time steps for observing f_n at $n = 1, \dots, N$. Further, we define a set of zero mean, mutually and serially uncorrelated random process $\{a_j\}$ and $\{b_j\}$ as finite sample analogues of the true $a(\omega)$ and $b(\omega)$ for the j -th spectral component. In particular, these processes are constant over n by covariance stationarity of f . Then, the finite sample analogue for the spectral representation theorem becomes [34]:

$$f_n = \mu_f + \sum_{j=1}^J [a_j \cos(\omega_j n) + b_j \sin(\omega_j n)] \quad (C7)$$

$$\mathbb{E}[a_j] = \mathbb{E}[b_j] = 0 \quad (C8)$$

$$\mathbb{E}[a_j a_{j'}] = \mathbb{E}[b_j b_{j'}] = \sigma^2 \delta(j - j') \quad (C9)$$

$$\mathbb{E}[a_j b_{j'}] = 0 \quad \forall j, j' \quad (C10)$$

$$\mu_f \equiv 0 \quad (C11)$$

The last line enforces a zero mean stochastic process and simplifies analysis without loss of generality and $\delta(\cdot)$ is the Dirac-delta function.

To illustrate, the first two moments are of the form:

$$\mathbb{E}[f_n] = \mu_f + \sum_{j=0}^J E[a_j] \cos(\omega_j n) + E[b_j] \sin(\omega_j n) = 0 \quad (\text{C12})$$

$$R(\nu) = \sum_j^J \sum_{j'}^J \sigma_j^2 \delta_{j,j'} [\cos(\omega_j n) \cos(\omega_{j'}(n + \nu)) + \sin(\omega_j n) \sin(\omega_{j'}(n + \nu))] \quad (\text{C13})$$

$$= \sigma^2 \sum_j^J p_j \cos(\omega_j \nu) \quad (\text{C14})$$

$$p_j \equiv \frac{\sigma_j^2}{\sigma^2} \equiv \frac{\sigma_j^2}{\sum_j \sigma_j^2} \quad (\text{C15})$$

We introduce process noise, w_n , into the formula for true f_n , and this establishes a commonality with state dynamics in Kalman filtering for a covariance stationary process:

$$f_n = \mu_f + \sum_{j=1}^J [a_j \cos(\omega_j(n-1)) + b_j \sin(\omega_j(n-1))] + w_n \quad (\text{C16})$$

In the absence of measurement noise and operating in the oversampling regime, an ordinary least squares (OLS) regression can be constructed by providing a collection of $J^{(B)}$ basis frequencies $\{\omega_j^{(B)}\}$, as in [34]. The OLS problem is constructed by separating the set of coefficients $\{\hat{\mu}_f, \hat{a}_1, \hat{b}_1, \dots, \hat{a}_J, \hat{b}_J\}$ and regressors $\{1, \cos(\omega_1(n-1)), \sin(\omega_1(n-1)), \dots, \cos(\omega_J^{(B)}(n-1)), \sin(\omega_J^{(B)}(n-1))\}$. For the specific particular choice of basis, $J^{(B)} = (N-1)/2$, (odd N) and $\omega_j^{(B)} \equiv 2\pi j/N$, we state the key result from [34] that the coefficient estimates are obtained as:

$$\hat{f}_n = \hat{\mu}_f + \sum_{j=1}^{J^{(B)}} [\hat{a}_j \cos(\omega_j^{(B)}(n-1)) + \hat{b}_j \sin(\omega_j^{(B)}(n-1))] \quad (\text{C17})$$

$$\hat{a}_j \equiv \frac{2}{N} \sum_{n'=1}^N \hat{f}_{n'} \cos(\omega_j^{(B)}(n'-1)) \quad (\text{C18})$$

$$\hat{b}_j \equiv \frac{2}{N} \sum_{n'=1}^N \hat{f}_{n'} \sin(\omega_j^{(B)}(n'-1)) \quad (\text{C19})$$

This choice of basis results in the number of regressors being the same as the length of the measurement record. Further, the term $(\hat{a}_j^2 + \hat{b}_j^2)$ is proportional to the total contribution of the j -th spectral component to the total sample variance of f , or in other words, the amplitude estimate for the power spectral density of true f .

Next, we depart from the OLS problem above by in several ways, firstly, by introducing measurement noise and secondly, by changing basis oscillators considered in the problem above. As in the main text, the linear measurement record is defined as:

$$y_n \equiv f_n + v_n \quad (\text{C20})$$

$$v_n \sim \mathcal{N}(0, R) \quad (\text{C21})$$

The link in GPR (periodic kernel) is direct and the link with LKFFB is made by setting $f_n \equiv H_n x_n$. In both frameworks, we incorporate the effect of measurement noise through the measurement noise variance, R , which has the effect of regularising the least squares estimation process discussed above.

1. Infinite Basis of Oscillators in a GPR Periodic Kernel

The departure from simple OLS plus measurement noise (above) to GPR (periodic kernel) arises from the fact that data is projected on an infinite basis of oscillators, namely, $J^{(B)} \rightarrow \infty$.

We follow the sketch of a proof provided in [42] to show that a sine squared exponential (periodic kernel) used in Gaussian Process Regression satisfies covariance function of trigonometric polynomials. Here, the index j labels an infinite comb of oscillators and m represents the higher order terms in the power reduction formulae in the last line of the definition below:

$$\omega_0^{(B)} \equiv \frac{\omega_j^{(B)}}{j}, j \in \{0, 1, \dots, J^{(B)}\} \quad (C22)$$

$$R(\nu) \equiv \sigma^2 \exp\left(-\frac{2 \sin^2\left(\frac{\omega_0^{(B)} \nu}{2}\right)}{l^2}\right) \quad (C23)$$

$$= \sigma^2 \exp\left(-\frac{1}{l^2}\right) \exp\left(\frac{\cos(\omega_0^{(B)} \nu)}{l^2}\right) \quad (C24)$$

$$= \sigma^2 \exp\left(-\frac{1}{l^2}\right) \sum_{m=0}^{M \rightarrow \infty} \frac{1}{m!} \frac{\cos^m(\omega_0^{(B)} \nu)}{l^{2m}} \quad (C25)$$

Next, we expand each cosine using power reduction formulae for odd and even powers respectively, and we re-group terms. For example, we expand the terms for $m = 0, 1, 2, 3, 4, 5 \dots$ as:

$$R(\nu) = \sigma^2 \exp\left(-\frac{1}{l^2}\right) \cos(\omega_0^{(B)} \nu) \left[\frac{2}{(2l^2)} \binom{1}{0} + \frac{2}{(2l^2)^3} \frac{1}{3!} \binom{3}{1} + \frac{2}{(2l^2)^5} \frac{1}{5!} \binom{5}{2} \dots \right] \quad (C26)$$

$$+ \sigma^2 \exp\left(-\frac{1}{l^2}\right) \cos(2\omega_0^{(B)} \nu) \left[\frac{2}{(2l^2)^2} \frac{1}{2!} \binom{2}{0} + \frac{2}{(2l^2)^4} \frac{1}{4!} \binom{4}{1} + \dots \right] \quad (C27)$$

$$+ \sigma^2 \exp\left(-\frac{1}{l^2}\right) \cos(3\omega_0^{(B)} \nu) \left[\frac{2}{(2l^2)^3} \frac{1}{3!} \binom{3}{0} + \frac{2}{(2l^2)^5} \frac{1}{5!} \binom{5}{1} \dots \right] \quad (C28)$$

$$+ \sigma^2 \exp\left(-\frac{1}{l^2}\right) \cos(4\omega_0^{(B)} \nu) \left[\frac{2}{(2l^2)^4} \frac{1}{4!} \binom{4}{0} + \dots \right] \quad (C29)$$

$$+ \sigma^2 \exp\left(-\frac{1}{l^2}\right) \cos(5\omega_0^{(B)} \nu) \left[\frac{2}{(2l^2)^5} \frac{1}{5!} \binom{5}{0} + \dots \right] \quad (C30)$$

⋮

$$+ \sigma^2 \exp\left(-\frac{1}{l^2}\right) \left[\frac{1}{(2l^2)^2} \frac{1}{2!} \binom{2}{1} + \frac{1}{(2l^2)^4} \frac{1}{4!} \binom{4}{2} + \dots \right] + \sigma^2 \exp\left(-\frac{1}{l^2}\right) \quad (C31)$$

In the expansion above, the vertical and horizontal dots represent contributions from $m > 5$ terms. The key message is that truncating m to a finite number of terms M will truncate j to represent a finite number of oscillators. For the example above, if the power reduction expansion indexed by m above was truncated to $M = 5$ terms, then the number of basis oscillators (number of rows) would also be truncated. We now summarise the amplitudes Eq. (C26) to Eq. (C30) in second term of $R(\nu)$ and Eq. (C31) corresponds to $p_{0,M}$ term below:

$$R(\nu) = \sigma^2 (p_{0,M} + \sum_{j=0}^{\infty} p_{j,M} \cos(j\omega_0^{(B)} \nu)) \quad (C32)$$

$$p_{j,M} \equiv \sigma^2 \exp\left(-\frac{1}{l^2}\right) \sum_{\beta=0}^{\beta=\beta_{j,m}^{MAX}} \frac{2}{(2l^2)^{(j+2\beta)}} \frac{1}{(j+2\beta)!} \binom{j+2\beta}{\beta} \quad (C33)$$

$$\beta \equiv 0, 1, \dots, \beta_{j,m}^{MAX} \quad (C34)$$

$$p_{0,M} = \exp\left(-\frac{1}{l^2}\right) \sum_{\alpha=0}^{\alpha=\alpha_m^{MAX}} \frac{1}{(2l^2)^{(2\alpha)}} \frac{1}{(2\alpha)!} \binom{2\alpha}{\alpha} \quad (C35)$$

$$\alpha \equiv 0, 1, \dots, \alpha_m^{MAX} \quad (C36)$$

By examining the cosine expansion, one sees that a truncation at $(M, J^{(B)})$ means our summarised formulae will require $\beta_{j,M}^{MAX} = \lfloor \frac{M-j}{2} \rfloor$ and $\alpha_M^{MAX} = \lfloor \frac{M}{2} \rfloor$ where $\lfloor \cdot \rfloor$ denotes the ceiling floor. If we truncate with $M \equiv J^{(B)}$ such that $\alpha_M^{MAX} = \lfloor \frac{J^{(B)}}{2} \rfloor$, $\beta_{j,M}^{MAX} = \lfloor \frac{J^{(B)}-j}{2} \rfloor$ and re-adjust the kernel for the zero-th frequency term, then we agree with final result in [42].

We compare the covariance function of the periodic kernel in Eq. (C32) with the covariance function of trigonometric polynomials in Eq. (C14). Here, $p_{j,M}$ for the periodic kernel are not identically specified in general to those under the spectral representation theorem, but otherwise retain a structure as a cosine basis where the correlations between two random variables in a sequence only depends on the separation between them. For a constant mean Gaussian process, the form of the periodic kernel allows the underlying process to satisfy covariance stationarity and appears to permit an interpretation via the spectral representation theorem.

2. Amplitude and Phase Extraction for Finite Oscillator Basis in LKFFB

In LKFFB, we depart from the simple OLS plus measurement noise problem considered earlier by specifying a fixed basis of oscillators at the physical Fourier resolution established by the measurement record. Using a specific state space model, we can track amplitudes and phases for each basis oscillator individually to enable forward prediction at any time-step of our choosing. The design of a fixed basis necessarily incorporates a priori assumptions about the extent to which a fast measurement action over-samples slowly drifting non-Markovian noise, that is, a (potentially incorrect) assumption about dephasing noise bandwidth.

The efficacy of the Liska Kalman Filter in our application assumes an appropriate choice of the ‘Kalman basis’ oscillators. The choice of basis can effect the forward prediction of state estimates. To illustrate, consider the choice of Basis A - C defined below. Basis A depicts a constant spacing above the Fourier resolution (e.g. $\omega_0^{(B)} \geq \frac{2\pi}{N_T \Delta t}$). Basis B introduces a minimum Fourier resolution and effectively creates an irregular spacing if one wishes to consider a basis frequency comb coarser than the experimentally established Fourier spacing over the course of the experiment. Basis C is identical to Basis B but allows a projection to a zero frequency component.

$$\text{Basis A: } \equiv \{0, \omega_0^{(B)}, 2\omega_0^{(B)} \dots J^{(B)}\omega_0^{(B)}\} \quad (\text{C37})$$

$$\text{Basis B: } \equiv \left\{ \frac{2\pi}{N\Delta t}, \frac{2\pi}{N\Delta t} + \omega_0^{(B)}, \dots, \frac{2\pi}{N\Delta t} + J^{(B)}\omega_0^{(B)} \right\} \quad (\text{C38})$$

$$\text{Basis C: } \equiv \left\{ 0, \frac{2\pi}{N\Delta t}, \frac{2\pi}{N\Delta t} + \omega_0^{(B)}, \dots, \frac{2\pi}{N\Delta t} + J^{(B)}\omega_0^{(B)} \right\} \quad (\text{C39})$$

While one can propagate LKFFB with zero gain, it may be advantageous for predictive control applications to generate predictions in one calculation rather than recursively. This means we sum contributions over all $j \in J^{(B)}$ oscillators and we reconstruct the signal for all future time values in one calculation, without having to propagate the filter recursively with zero gain. The interpretation of the predicted signal, \hat{f}_n , requires an additional (but time-constant) phase correction term ψ_C that arises as a byproduct of the computational basis (i.e. Basis A, B or C). The phase correction term corrects for a gradual mis-alignment between Fourier and computational grids which occurs if one specifies a non-regular spacing inherent in Basis B or C. Let n_C denote the time-step at which instantaneous amplitudes $\|\hat{x}_{n_C}^j\|$ and instantaneous phase $\theta_{\hat{x}_{n_C}^j}$ is extracted for the oscillator represented by the j -th state space resonator, x_n^j , where super-script j denotes an oscillator of frequency $\omega_j^{(B)} \equiv j\omega_0^{(B)}$ (not a power):

$$\hat{f} = \sum_{j=0}^{J^{(B)}} \|\hat{x}_{n_C}^j\| \cos(m\Delta t\omega_j^{(B)} + \theta_{\hat{x}_{n_C}^j} + \psi_C), \quad (\text{C40})$$

$$n_C \in N_T, \quad m \in N_P$$

$$\psi_C \equiv \begin{cases} 0, & (\text{Basis A}) \\ \equiv \frac{2\pi}{\omega_0^{(B)}}(\omega_0^{(B)} - \frac{2\pi}{N\Delta t}), & (\text{Basis B or C}) \end{cases} \quad (\text{C41})$$

The output predictions from calculating a harmonic sum using learned instantaneous amplitudes, phases and the LKFFB Basis A-C agree with zero-gain predictions if ψ_C is specified as above. The calculation of ψ_C is determined entirely by the choice of computational and experimental sampling procedures, and assumes no information about true dephasing.

Next, we define an analytical ratio to define the optimal training time, n_C , at which LKFFB predictions should commence, irrespective of whether the prediction procedure is recursively propagating the Kalman Filter with zero gain, or by calculating a harmonic sum for all prediction points in one go.

$$n_C \equiv \frac{1}{\Delta t\omega_0^{(B)}} = \frac{f_s}{\omega_0^{(B)}} \quad (\text{C42})$$

Consider an arbitrarily chosen training period, $N_T \neq n_C$. For f_s fixed, our choice of $N_T > n_C$ means we are achieving a Fourier resolution which exceeds the resolution of the LKFFB basis. Now consider $N_T < n_C$. This means that we've extracted information prematurely, and we have not waited long enough to project on the smallest basis frequency, namely, $\omega_0^{(B)}$. In the case where data is perfectly projected on our basis, this has no impact. For imperfect learning, we see that instantaneous amplitude and phase information slowly degrades for $N_T > n_C$; and trajectories for the smallest basis frequency have not stabilised for $N_T < n_C$.

Of these choices, Basis A for $\omega_0^{(B)} \equiv \frac{2\pi}{N_T \Delta t}$ is expected to yield best performance, at the expense of computational load, and this is confirmed in numerical experiments. All results in this manuscript are reported for Basis A with $N_T \equiv \frac{1}{\Delta t \omega_0^{(B)}} = \frac{f_s}{\omega_0^{(B)}}$.

3. Equivalent Spectral Representation of f in LKFFB and GPR Periodic Kernel

In this section, we consider the structural similarities between LKFFB and GPR with a periodic kernel. We show that the LKFFB has an analogous structure to a stack of stochastic processes on a circle [38], and in moving from discrete to continuous time, we recover a covariance function that has the same structure if the periodic kernel was truncated to a finite basis of oscillators, $J^{(B)}$. For zero mean, Gaussian random variables, covariance stationarity is established, completing the link between LKFFB and the periodic kernel. For the case $\Gamma_n w_n \rightarrow w_n$ in LKFFB, stacked Kalman resonators as an approximation to infinite oscillators in a periodic kernel is documented in [42].

At time step n , the posterior Kalman state at $n-1$ acts as the initial state at n , such that $\nu = \Delta t$ for a small Δt such that a linearised trajectory is approximately true for each basis frequency. We show this using the following correlation relations and a Gaussian assumption for process noise, where $n, m \in N$ are indices for time steps and $j = 0, 1, \dots, J^{(B)}$ indexes the set of basis oscillators:

$$\mathbb{E}[w_n] = 0 \quad \forall j \in J^{(B)}, \quad n \in N \quad (\text{C43})$$

$$\mathbb{E}[w_n, w_m] = \sigma^2 \delta(n - m) \quad n, m \in N \quad (\text{C44})$$

$$\mathbb{E}[A_0^j] = \mathbb{E}[B_0^{j'}] = 0, \quad \forall j, j' \in J^{(B)} \quad (\text{C45})$$

$$\mathbb{E}[A_n^j B_m^{j'}] = 0, \quad \forall j, j' \in J^{(B)}, \quad n, m \in N \quad (\text{C46})$$

$$\mathbb{E}[A_n^j A_m^{j'}] = \mathbb{E}[B_n^j B_m^{j'}] = \sigma_j^2 \delta(n - m) \delta(j - j'), \quad \forall j, j' \in J^{(B)}, \quad n, m \in N \quad (\text{C47})$$

$$\mathbb{E}[w_n A_m^j] = \mathbb{E}[w_n B_m^{j'}] \equiv 0 \quad \forall j, j' \in J^{(B)}, \quad n, m \in N \quad (\text{C48})$$

Consider a j -th state space resonator, x_n^j , in the LKFFB, where super-script j denotes an oscillator (not a power) and we obtain:

$$\Theta(j\omega_0^{(B)} \Delta t) = \begin{bmatrix} \cos(j\omega_0^{(B)} \Delta t) & -\sin(j\omega_0^{(B)} \Delta t) \\ \sin(j\omega_0^{(B)} \Delta t) & \cos(j\omega_0^{(B)} \Delta t) \end{bmatrix} \quad (\text{C49})$$

$$x_n^j \equiv \begin{bmatrix} A_n^j \\ B_n^j \end{bmatrix} = \Theta(j\omega_0^{(B)} \Delta t) \left[\hat{\mathcal{I}} + \frac{w_{n-1}}{\sqrt{A_{n-1}^{j^2} + B_{n-1}^{j^2}}} \right] \begin{bmatrix} A_{n-1}^j \\ B_{n-1}^j \end{bmatrix} \quad (\text{C50})$$

$$(\text{C51})$$

$$\implies \mathbb{E}[x_n^j] = 0 \quad (\text{C52})$$

$$\implies \mathbb{E}[x_n^j x_m^j]^T = \Theta(j\omega_0^{(B)} \Delta t) \mathbb{E} \left[\begin{bmatrix} A_{n-1}^j A_{m-1}^j & A_{n-1}^j B_{m-1}^j \\ B_{n-1}^j A_{m-1}^j & B_{n-1}^j B_{m-1}^j \end{bmatrix} \right] \Theta(j\omega_0^{(B)} \Delta t)^T \quad (\text{C53})$$

$$+ \Theta(j\omega_0^{(B)} \Delta t) \left[\frac{w_{n-1}}{\sqrt{A_{n-1}^j{}^2 + B_{n-1}^j{}^2}} + \frac{w_{m-1}}{\sqrt{A_{m-1}^j{}^2 + B_{m-1}^j{}^2}} \right] \begin{bmatrix} A_{n-1}^j A_{m-1}^j & A_{n-1}^j B_{m-1}^j \\ B_{n-1}^j A_{m-1}^j & B_{n-1}^j B_{m-1}^j \end{bmatrix} \Theta(j\omega_0^{(B)} \Delta t)^T \quad (\text{C54})$$

$$+ \Theta(j\omega_0^{(B)} \Delta t) \left[\frac{w_{n-1} w_{m-1}}{\sqrt{A_{n-1}^j{}^2 + B_{n-1}^j{}^2} \sqrt{A_{m-1}^j{}^2 + B_{m-1}^j{}^2}} \right] \begin{bmatrix} A_{n-1}^j A_{m-1}^j & A_{n-1}^j B_{m-1}^j \\ B_{n-1}^j A_{m-1}^j & B_{n-1}^j B_{m-1}^j \end{bmatrix} \Theta(j\omega_0^{(B)} \Delta t)^T \quad (\text{C55})$$

$$= \sigma_j^2 \delta(n-m) \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (\text{C56})$$

The cross correlation terms disappear under the temporal correlation functions so defined, namely, if assume $n \geq m$, then states A_{m-1}^j, B_{m-1}^j at $m-1$ at most have a w_{n-2} term (for the case $n = m$) and cannot be correlated with a future noise term w_{n-1} .

The dynamical trajectory in LKFFB is linearised for small Δt . The linearisation is an approximation to a true, continuous time deterministic trajectory defining a stochastic process on a circle.

We briefly visit this continuous time trajectory to specify the link between LKFFB and GPR (periodic kernel). Let t denote the continuous time deterministic dynamics for random initial state given by a_0^j, b_0^j , where super-script j denotes an oscillator with frequency $\omega_j \equiv j\omega_0^{(B)}$ (not a power):

$$\mathbb{E}[a_0^j] = \mathbb{E}[b_0^{j'}] = 0, \quad \forall j, j' \in J^{(B)} \quad (\text{C57})$$

$$\mathbb{E}[a_0^j b_0^{j'}] = 0, \quad \forall j, j' \in J^{(B)} \quad (\text{C58})$$

$$\mathbb{E}[a_0^j a_0^{j'}] = \mathbb{E}[b_0^j b_0^{j'}] = \sigma_j^2 \delta(j - j'), \quad \forall j, j' \in J^{(B)} \quad (\text{C59})$$

$$x^j(t) \equiv \begin{bmatrix} \cos(\omega_j t) & -\sin(\omega_j t) \\ \sin(\omega_j t) & \cos(\omega_j t) \end{bmatrix} \begin{bmatrix} a_0^j \\ b_0^j \end{bmatrix} \quad (\text{C60})$$

$$E[x^j(t)] = 0 \quad (\text{C61})$$

$$E[x^j(t) x^j(t')^T] = \begin{bmatrix} \cos(\omega_j t') & -\sin(\omega_j t') \\ \sin(\omega_j t') & \cos(\omega_j t') \end{bmatrix} \begin{bmatrix} a_0^j \\ b_0^j \end{bmatrix} \begin{bmatrix} a_0^j & b_0^j \end{bmatrix} \begin{bmatrix} \cos(\omega_j t) & -\sin(\omega_j t) \\ \sin(\omega_j t) & \cos(\omega_j t) \end{bmatrix} \quad (\text{C62})$$

$$= \sigma_j^2 \begin{bmatrix} \cos(\omega_j \nu) & 0 \\ 0 & \cos(\omega_j \nu) \end{bmatrix}, \quad \nu \equiv |t' - t| \quad (\text{C63})$$

We see that the initial state variables, a_0^j, b_0^j , must be zero mean, independent and identically distributed variables for each j such that $x^j(t)$ is covariance stationary. If a_0^j, b_0^j are Gaussian, then the joint distribution, $x^j(t)$, remains Gaussian under the linear operations above. Hence, the continuous time limit of the dynamics in LKFFB for $J^{(B)}$ independent substates, $x^j(t)$, describe a process with the same first and second moments for a periodic kernel truncated at $J^{(B)}$. For Gaussian processes, this results in an approximate equivalent representation of LKFFB for $J^{(B)}$ stacked resonators with an expansion of the periodic kernel truncated at $J^{(B)}$.

While the formalism of LKFFB shares a common structure with GPR (periodic kernel) in a particular limit, the physical interpretation of A_n^j, B_n^j is that these are components of the Hilbert transform of the original signal [29]. This gives us the ability to track and extract instantaneous amplitude and phase associated with each basis oscillator in LKFFB. In contrast, the coefficients of the periodic kernel are always contingent on the arbitrary truncation of the infinite basis, as seen in Eqs. (C32), (C33) and (C35). Hence, tracking (or extracting) amplitudes and phases for individual oscillators does not seem appropriate for the periodic kernel, as these values would change depending on the arbitrary choice of a truncation point.

- [1] J. J. J. Groen, R. Paap, and F. Ravazzolo, "Real-time inflation forecasting in a changing world," *Journal of Business & Economic Statistics* **31**, 29 (2013).
- [2] Y. Dong, Y. Li, M. Xiao, and M. Lai, "Unscented Kalman filter for time varying spectral analysis of earthquake ground motions," *Applied Mathematical Modelling* **33**, 398 (2009).
- [3] J. Ko and D. Fox, "GP BayesFilters: Bayesian filtering using Gaussian process prediction and observation models," *Autonomous Robots* **27**, 75 (2009).
- [4] A. C. Harvey, *Forecasting, structural time series models and the Kalman filter* (Cambridge University Press, Cambridge, United Kingdom, 1990).
- [5] C. Cheng, A. Sa-Ngasoongsong, O. Beyca, T. Le, H. Yang, Z. Kong, and S. T. Bukkapatnam, "Time series forecasting for nonlinear and non-stationary processes: a review and comparative study," *IIE Transactions* **47**, 1053 (2015).
- [6] J. D. Garcia and G. C. Amaral, "An optimal polarization tracking algorithm for Lithium-Niobate-based polarization controllers," in *Sensor Array and Multichannel Signal Processing Workshop (SAM)* (Rio de Janeiro, 2016) pp. 1–5.
- [7] F. R. Bach and M. I. Jordan, "Learning graphical models for stationary time series," *IEEE transactions on signal processing* **52**, 2189 (2004).
- [8] S. Tatinati and K. C. Veluvolu, "A hybrid approach for short-term forecasting of wind speed," *The Scientific World Journal* **2013**, 548370 (2013).
- [9] J. Hall, C. E. Rasmussen, and J. Maciejowski, "Reinforcement learning with reference tracking control in continuous state spaces," in *Decision and Control and European Control Conference (CDC-ECC)* (Orlando, 2011) pp. 6019–6024.
- [10] F. Hamilton, T. Berry, and T. Sauer, "Ensemble Kalman filtering without a model," *Physical Review X* **6**, 011021 (2016).
- [11] J. V. Candy, *Bayesian signal processing: classical, modern, and particle filtering methods*, Vol. 54 (John Wiley & Sons, Hoboken, New Jersey, 2016).
- [12] B. Boashash, "Estimating and interpreting the instantaneous frequency of a signal. II. Algorithms and applications," *Proceedings of the IEEE* **80**, 540 (1992).
- [13] L. Ji and Z. Tie, "On gradient descent algorithm for generalized phase retrieval problem," in *Proceedings of IEEE 13th International Conference on Signal Processing (ICSP)* (IEEE, China, 2016) pp. 320–325.
- [14] G. Struchalin, I. Pogorelov, S. Straupe, K. Kravtsov, I. Radchenko, and S. Kulik, "Experimental adaptive quantum tomography of two-qubit states," *Physical Review A* **93**, 012103 (2016).
- [15] A. Sergeevich, A. Chandran, J. Combes, S. D. Bartlett, and H. M. Wiseman, "Characterization of a qubit Hamiltonian using adaptive measurements in a fixed basis," *Physical Review A* **84**, 052315 (2011).
- [16] D. Mahler, L. A. Rozema, A. Darabi, C. Ferrie, R. Blume-Kohout, and A. Steinberg, "Adaptive quantum state tomography improves accuracy quadratically," *Physical Review Letters* **111**, 183601 (2013).
- [17] M. P. Stenberg, O. Köhn, and F. K. Wilhelm, "Characterization of decohering quantum systems: Machine learning approach," *Physical Review A* **93**, 012122 (2016).
- [18] A. Shabani, R. Kosut, M. Mohseni, H. Rabitz, M. Broome, M. Almeida, A. Fedrizzi, and A. White, "Efficient measurement of quantum dynamics via compressive sensing," *Physical Review Letters* **106**, 100401 (2011).
- [19] Z. Shen, W. X. Wang, Y. Fan, Z. Di, and Y.-C. Lai, "Reconstructing propagation networks with natural diversity and identifying hidden sources," *Nature Communications* **5** (2014), 10.1038/ncomms5323.
- [20] L. E. de Clercq, R. Oswald, C. Flühmann, B. Keitch, D. Kienzler, H.-Y. Lo, M. Marinelli, D. Nadlinger, V. Negnevitsky, and J. P. Home, "Estimation of a general time-dependent Hamiltonian for a single qubit," *Nature Communications* **7** (2016), 10.1038/ncomms11218.
- [21] D. Tan, S. Weber, I. Siddiqi, K. Moelmer, and K. Murch, "Prediction and retrodiction for a continuously monitored superconducting qubit," *Physical Review Letters* **114**, 090403 (2015).
- [22] Y. Huang and J. E. Moore, "Neural network representation of tensor network and chiral states," arXiv:1701.06246.
- [23] C. Bonato, M. S. Blok, H. T. Dinani, D. W. Berry, M. L. Markham, D. J. Twitchen, and R. Hanson, "Optimized quantum sensing with a single electron spin using real-time adaptive measurements," *Nature Nanotechnology* **11**, 247 (2016).
- [24] N. Wiebe, C. Granade, A. Kapoor, and K. M. Svore, "Bayesian inference via rejection filtering," arXiv:1511.06458.
- [25] M. D. Shulman, S. P. Harvey, J. M. Nichol, S. D. Bartlett, A. C. Doherty, V. Umansky, and A. Yacoby, "Suppressing qubit dephasing using real-time Hamiltonian estimation," *Nature Communications* **5** (2014), 10.1038/ncomms6156.
- [26] C. Granade, J. Combes, and D. Cory, "Practical Bayesian tomography," *New Journal of Physics* **18**, 033024 (2016).
- [27] P. E. Jacob, S. M. M. Alavi, A. Mahdi, S. J. Payne, and D. A. Howey, "Bayesian inference in non-Markovian state-space models with applications to battery fractional-order systems," *IEEE Transactions on Control Systems Technology* **26**, 497 (2017).
- [28] S. Mavadia, V. Frey, S. D. Jarrah Sastrawan, and M. J. Biercuk, "Prediction and real-time compensation of qubit decoherence via machine learning," *Nature Communications* **8** (2017), 10.1038/ncomms14106.
- [29] J. Liška and E. Janeček, "Time-frequency representation of instantaneous frequency using a Kalman filter," in *Robotics Automation and Control* (InTech, Vienna, 2008) pp. 28–38.
- [30] C. Ferrie, C. E. Granade, and D. G. Cory, "How to best sample a periodic probability distribution or on the accuracy of Hamiltonian finding strategies," *Quantum Information Processing* **12**, 611 (2013).
- [31] M. S. Grewal and A. P. Andrews, *Kalman Filtering: Theory and Practice Using MATLAB*, 2nd ed. (John Wiley & Sons, Hoboken, New Jersey, 2001).
- [32] S.-M. Moon, D. G. Cole, and R. L. Clark, "Real-time implementation of adaptive feedback and feedforward gen-

- eralized predictive control algorithm,” *Journal of Sound and Vibration* **294**, 82 (2006).
- [33] I. D. Landau, R. Lozano, M. M’Saad, and A. Karimi, *Adaptive control: Algorithms, analysis and applications*, Vol. 51 (Springer, Berlin, 1998).
 - [34] J. D. Hamilton, *Time Series Analysis*, Vol. 2 (Princeton University Press, Princeton, New Jersey, 1994).
 - [35] P. J. Brockwell and R. A. Davis, *Introduction to Time Series and Forecasting* (Springer-Verlag, New York, 1996).
 - [36] M. Salzmann, P. Teunissen, and M. Sideris, “Detection and modelling of coloured noise for Kalman filter applications,” in *Kinematic Systems in Geodesy, Surveying, and Remote Sensing*, Vol. 107 (Springer-Verlag, New York, 1991) pp. 251–260.
 - [37] B. Wahlberg, “Estimation of autoregressive moving-average models via high-order autoregressive approximations,” *Journal of Time Series Analysis* **10**, 283 (1989).
 - [38] S. Karlin and H. Taylor, *A First Course in Stochastic Processes* (Academic Press Inc, New York, 1975).
 - [39] R. Karlsson and F. Gustafsson, *Filtering and estimation for quantized sensor information*, Tech. Rep. LiTH-ISY-R-2674 (Linköping University, 2005).
 - [40] B. Widrow, I. Kollar, and M. C. Liu, “Statistical theory of quantization,” *IEEE Transactions on Instrumentation and Measurement* **45**, 353 (1996).
 - [41] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)* (MIT Press, Cambridge, 2005).
 - [42] A. Solin and S. Särkkä, “Explicit link between periodic covariance functions and state space models,” in *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics*, Proceedings of Machine Learning Research, Vol. 33, edited by S. Kaski and J. Corander (PMLR, Reykjavik, 2014) pp. 904–912.
 - [43] F. Tobar, T. D. Bui, and R. E. Turner, “Learning stationary time series using Gaussian processes with non-parametric kernels,” in *Advances in Neural Information Processing Systems*, Vol. 28 (Curran Associates, New York, 2015) pp. 3501–3509.
 - [44] S. Roberts, M. Osborne, M. Ebden, S. Reece, N. Gibson, and S. Aigrain, “Gaussian processes for time-series modelling,” *Phil. Trans. R. Soc. A* **371**, 20110550 (2013).
 - [45] M. L. Stein, *Interpolation of Spatial Data: Some Theory for Kriging* (Springer Science & Business Media, 1999).
 - [46] A. Soare, H. Ball, D. Hayes, X. Zhen, M. Jarratt, J. Sastrowan, H. Uys, and M. Biercuk, “Experimental bath engineering for quantitative studies of quantum control,” *Physical Review A* **89**, 042329 (2014).
 - [47] GPy, “GPy: A Gaussian process framework in Python,” <http://github.com/SheffieldML/GPy> (2012).
 - [48] S. Arlot and P. Massart, “Data-driven calibration of penalties for least-squares regression,” *Journal of Machine Learning Research*, **10**, 245 (2009).
 - [49] K. Vu, J. C. Snyder, L. Li, M. Rupp, B. F. Chen, T. Kheif, K.-R. Müller, and K. Burke, “Understanding kernel ridge regression: Common behaviors from simple functions to density functionals,” *International Journal of Quantum Chemistry* **115**, 1115 (2015).
 - [50] P. Abbeel, A. Coates, M. Montemerlo, A. Y. Ng, and S. Thrun, “Discriminative training of Kalman filters,” in *Robotics: Science and Systems* (MIT Press, Cambridge, 2005) pp. 289–296.
 - [51] A. Robertson and C. Grenade, (unpublished).
 - [52] A. Wilson and R. Adams, “Gaussian process kernels for pattern discovery and extrapolation,” in *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, Vol. 28 (Journal of Machine Learning Research, Atlanta, 2013) pp. 1067–1075.
 - [53] J. Quiñero Candela, C. E. Rasmussen, A. R. Figueiras-Vidal, and M. Lázaro-Gredilla, “Sparse spectrum Gaussian process regression,” *Journal of Machine Learning Research* **11**, 1865 (2010).
 - [54] A. Gelb, *Applied Optimal Estimation* (MIT Press, Cambridge, 1974).
 - [55] M. West and J. Harrison, *Bayesian Forecasting and Dynamic Models* (Springer-Verlag, New York, 1996).