

# Finding Birkhoff Averages via Adaptive Filtering

M. Ruth<sup>1, a)</sup> and D. Bindel<sup>2</sup>

<sup>1)</sup>*Center for Applied Mathematics, Cornell University, Ithaca, NY*

<sup>2)</sup>*Department of Computer Science, Cornell University, Ithaca, NY*

(Dated: 29 March 2024)

In many applications, one is interested in classifying trajectories of Hamiltonian systems as invariant tori, islands, or chaos. The convergence rate of ergodic Birkhoff averages can be used to categorize these regions, but many iterations of the return map are needed to implement this directly. Recently, it has been shown that a weighted Birkhoff average can be used to accelerate the convergence, resulting in a useful method for categorizing trajectories.

In this paper, we show how a modified version the reduced rank extrapolation method (named Birkhoff RRE) can also be used to find optimal weights for the weighted average with a single linear least-squares solve. Using these, we classify trajectories with fewer iterations of the map than the standard weighted Birkhoff average. Furthermore, for the islands and invariant circles, a subsequent eigenvalue problem gives the number of islands and the rotation number. Using these numbers, we find Fourier parameterizations of invariant circles and islands. We show examples of Birkhoff RRE on the standard map and on magnetic field line dynamics.

## I. INTRODUCTION

Invariant tori are ubiquitous structures in symplectic maps and Hamiltonian dynamics. Examples of invariant tori include periodic orbits of the pendulum, invariant circles in the standard map, halo orbits in astrodynamics<sup>1</sup>, and nested flux surfaces in magnetic confinement devices<sup>2</sup>. Such orbits are known to be stable to perturbation due to the KAM theorem<sup>3</sup>. However, numerically identifying orbits from trajectories is often challenging, due to both the difficulties of high-dimensional geometry and the problems of small denominators.

One standard numerical method to find invariant tori is the parameterization method<sup>4</sup>. This method is based on the conjugacy relation defining invariant tori, and can be accelerated using the fast Fourier transform. This allows for highly accurate computations of invariant tori, which can be proven to exist by a numerical variant of the KAM theorem<sup>5</sup>. However, one of the main drawbacks of the parameterization method is that it needs an initial guess. In the case of 1D and 1.5D Hamiltonian systems, a manual initial guess can be found relatively easily with a Poincaré plot of a trajectory. Once one solution is found, continuation<sup>1</sup> can be used to find more solutions. However, in the cases of higher dimensional systems or islands, initial guesses are significantly harder to find. The initial guess issue is faced by other methods relying on iterations on torus parameterizations, such as the flux minimizing surfaces<sup>6,7</sup>.

Additionally, initial guesses of the rotation number are also difficult. Simple methods for finding the rotation are available for invariant circles where there is a natural point to wind about. In such cases, one can find the rotation number via classical limits<sup>8</sup> or more accurately using

the weighted Birkhoff average<sup>9–11</sup> or Richardson-like extrapolation algorithms<sup>12,13</sup>. Once the rotation number is found, then it is straightforward to find a parameterization of the orbit<sup>14</sup>. Without good guesses at the winding structure, the typical solution is to look for peaks in the Fourier spectrum of the signal or to use another frequency-based method<sup>15</sup>. Unfortunately, these peaks are only as accurate as the discretization resolution of the spectrum. This issue is again made more complicated in higher dimensions, where winding is a less well-defined concept.

An additional issue that we will discuss is that of orbit classification. Before an invariant torus can be fit to a trajectory, we must first be confident that trajectory is, in fact, a torus. Classical methods for this typically rely on the Lyapunov exponent, but this can be quite slow to converge. More recently, the rapid convergence of the weighted Birkhoff average<sup>9–11,16–18</sup> (WBA) has been shown to be capable distinguishing chaotic from non-chaotic trajectories.

Symplectic maps are often computed by evolving Hamiltonian dynamics by numerical integration, and this can present its own problems. One problem is noise: reliable methods must be robust to potentially non-symplectic errors in the symplectic map. This is particularly relevant when symplectic integrators are not available for a given application. We note that for the parameterization method, this is alleviated by an overdetermined formulation of the method<sup>19</sup>. A second issue is the cost of evaluating the map. Both the parameterization method using the Fast Fourier Transform and WBA can be performed in nearly linear time in the number of samples, so the dominant cost is typically dominated by evolving Hamiltonian trajectories. When classifying many trajectories, any reduction in the number of evaluations is very useful. Structured symplectic interpolants such as the HenonNet<sup>20</sup>, SympNet<sup>21</sup>, or a Gaussian Process approach<sup>22</sup> can alleviate both issues of numerical

<sup>a)</sup>mer335@cornell.edu

integration. However, symplectic interpolants are necessarily nonlinear, and may require more data than simply finding the invariant circle. Linear interpolants can similarly be used, but require high accuracy, and the evaluation time of constructing the interpolant may still represent the dominant cost of any algorithm.

In this paper, we classify and parameterize invariant tori from single trajectories efficiently in the number of map evaluations. For the classification step, we will rely on a variant of the reduced rank extrapolation method<sup>23</sup>, which we will refer to as Birkhoff RRE. Birkhoff RRE works by attempting to find a filter (or linear model) for measurements on a trajectory using only the time-series information. Because Birkhoff RRE only depends on a single trajectory, it does not require any initial guesses or continuation of invariant tori. Additionally, Birkhoff RRE is written as a linear least-squares problem, meaning its implementation is straightforward and there is a residual indicating the fit of the linear model. We prove that when the trajectory is on an invariant torus, the residual of Birkhoff RRE approaches zero as rapidly as the weighted Birkhoff average, allowing for the same classification ability. Additionally, we show experimentally that RRE converges significantly faster than WBA on a set of trajectories of the standard map, with a large majority of trajectories being classified to highly accurate residuals below  $10^{-11}$  in fewer than 1000 iterations of the map.

For the parameterization step, we show experimentally that the frequencies filtered by the Birkhoff RRE filter are multiples of the rotation number. Using those frequencies, we numerically identify both the number of islands and rotation numbers for trajectories in 2D. Once this information is known, we project the signal back onto the corresponding Fourier modes, giving a parameterization of the invariant circle or island. This process is similar to the filter diagonalization method<sup>24</sup>.

We introduce necessary background in Sec. II, including a review of the weighted Birkhoff average. In Sec. III, we build upon WBA to introduce the Birkhoff RRE algorithm, stating the relevant convergence theorems. Then, in Sec. IV, we show two examples of the method. In the first example, we examine the convergence of Birkhoff RRE on the standard map, confirming the predicted rates from Sec. III. Then, in the second example, we show how the method can be applied to an example on a symplectic map obtained from a toroidal plasma confinement device known as a stellarator, showing how the method can be used in a real-world situation. Finally, we conclude in Sec. V.

## II. BACKGROUND

Let  $\mathbf{F} : X \rightarrow X$  be a map for a discrete-time dynamical system, where  $X$  is a  $C^M$  manifold. A  $d$ -dimensional *invariant torus* of  $\mathbf{F}$  is a function  $\mathbf{z} : \mathbb{T}^d \rightarrow X$  that

satisfies the conjugacy

$$\mathbf{F} \circ \mathbf{z} - \mathbf{z} \circ R_\omega = 0,$$

where  $R_\omega : \mathbb{T}^d \rightarrow \mathbb{T}^d$  is the rotation map  $R_\omega(\boldsymbol{\theta}) = \boldsymbol{\theta} + \boldsymbol{\omega}$  with  $\boldsymbol{\omega} \in \mathbb{T}^d$ . The invariant torus has an associated invariant measure  $\mu$  on its image  $X_0 \subseteq X$ . We assume that the invariant tori are smooth, i.e.  $\mathbf{z} \in C^M$  for some positive integer  $M$  or  $\mathbf{z} \in C^\infty$ . We refer to  $d = 1$  invariant tori as *invariant circles*.

We require the rotation vector  $\boldsymbol{\omega}$  satisfy the  $(c, \nu)$  Diophantine condition

$$|\boldsymbol{\omega} \cdot \mathbf{n} - m| \geq \frac{c}{\|\mathbf{n}\|^\nu} \quad \text{for all } \mathbf{n} \in \mathbb{Z}^d \setminus \{0\}, m \in \mathbb{Z}.$$

For  $d = 1$ , the Diophantine condition tells us the  $\omega \in \mathbb{R}$  is “sufficiently irrational.” In higher dimensions, the condition gives a measure of irrational independence, requiring that the rotation numbers also be far from being rational multiples of the others. We note that the Diophantine condition is satisfied by almost all rotation numbers under the Lebesgue measure<sup>3</sup>. So, almost all invariant tori are Diophantine for nested regions with shear.

We define an *island* to be a set of  $p \geq 1$  tori  $\mathbf{z}^{(j)} : \mathbb{T}^d \rightarrow X$  with  $1 \leq j \leq p$  such that for some  $\boldsymbol{\omega}$  the following conjugacy is satisfied:

$$\mathbf{F} \circ \mathbf{z}^{(j)} = \begin{cases} \mathbf{z}^{(j+1)}, & 1 \leq j < p, \\ \mathbf{z}^{(1)} \circ R_\omega, & j = p. \end{cases} \quad (1)$$

A direct result of the above definition is that each  $\mathbf{z}^{(j)}$  is an invariant torus of the map  $\mathbf{F}^p$  (and hence a  $p = 1$  island is equivalent to an invariant torus). Another special case of an island is a periodic orbit, where each  $\mathbf{z}^{(j)}$  is a constant function. If a trajectory is on an invariant torus or island, we call it *integrable*.

We note that the definitions of invariant tori and islands above do not require any special properties of the map. However, they are both commonly found in the case of *symplectic maps*, which preserve some symplectic 2-form under the pushforward of the map. In the case of symplectic maps, we refer to everything that is not integrable as *chaotic*. We note that this is a heuristic definition of chaos, as the theorems herein only guarantee convergence rates of certain methods for integrable trajectories. The converse (i.e. chaotic trajectories converge slowly for some definition of chaos) is still an open problem.

In Fig. 1, we plot the phase portrait of the Chirikov standard map  $F : (x_t, y_t) \mapsto (x_{t+1}, y_{t+1})$  on  $X = \mathbb{T} \times \mathbb{R}$  where

$$\begin{aligned} x_{t+1} &= x_t + y_{t+1} \pmod{1}, \\ y_{t+1} &= y_t - \frac{k_{\text{sm}}}{2\pi} \sin(2\pi x_t), \end{aligned} \quad (2)$$

and  $k_{\text{sm}} = 0.7$ . The map has invariant circles (e.g. nested about  $(0, 0)$ ) and the blue-to-green gradient of trajectories in the center), islands (e.g. the  $p = 2$  island chain

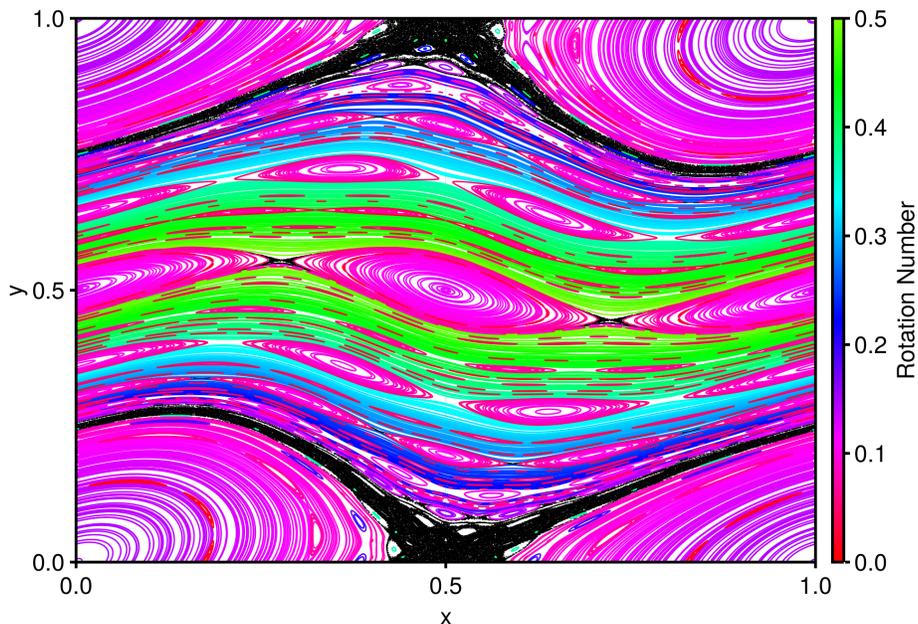


FIG. 1. Phase portrait of the standard map (2). Parameterizations of the invariant circles and islands are obtained via the methodology in Sec. III, using Algorithm 2 with  $\epsilon = 0$ ,  $\gamma = 3$ ,  $\delta = 10^{-10}$ ,  $K_{\text{init}} = 50$ ,  $K_{\text{max}} = 600$ , and  $\Delta K = 50$ . Invariant circles and islands are colored according to their rotation numbers, while trajectories classified as chaotic are plotted in black.

centered at  $(0.0, 0.5)$  and  $(0.5, 0.5)$ , and chaos in black (e.g. about  $(0.5, 0)$ ). The invariant circles and islands are both colored by the rotation number  $\omega$ , found using the methods in Sec. III. We note that  $\omega$  has a unique representation in  $[0.0, 0.5]$  due both to the fact that  $\omega \in \mathbb{T}$  and the freedom to take  $\theta \rightarrow -\theta$  in the parameterization.

To categorize orbits, we consider the problem of finding ergodic averages. Let  $\tilde{\mathbf{h}} : X \rightarrow \mathbb{R}^D$  be an observable function on our state space. The *Birkhoff average* of  $\tilde{\mathbf{h}}$  is defined as the limit of finite time averages:

$$\mathcal{B}[\tilde{\mathbf{h}}](\mathbf{x}) = \lim_{\bar{K} \rightarrow \infty} \mathcal{B}_{\bar{K}}[\tilde{\mathbf{h}}](\mathbf{x})$$

where

$$\mathcal{B}_{\bar{K}}[\tilde{\mathbf{h}}] = \frac{1}{\bar{K}} \sum_{k=0}^{\bar{K}-1} (\tilde{\mathbf{h}} \circ \mathbf{F}^k)(\mathbf{x}).$$

For an initial point  $\mathbf{x}$ , we let  $X_0$  be the ergodic component

$$X_0 = \{\mathbf{x}' \in X \mid \forall f \in C^b(X), \lim_{T \rightarrow \infty} \mathcal{B}_T[f](\mathbf{x}) - \mathcal{B}_T[f](\mathbf{x}') = 0\}$$

where  $C^b(X)$  is the set of continuous bounded functions on  $X$ . We note that the ergodic component  $X_0$  is identical to our previous definition for circles with irrational rotation numbers  $\omega$ . Then, for almost all  $\mathbf{x}$ , the Birkhoff average converges to an average over a unique invariant measure  $\mu$  on  $X_0$ <sup>25</sup>

$$\mathcal{B}[\tilde{\mathbf{h}}](\mathbf{x}) = \int_{X_0} \tilde{\mathbf{h}}(\mathbf{x}) d\mu.$$

Additionally, when  $\tilde{\mathbf{h}} \in C^M$  for  $M > \nu + d$ , one can show that the partial averages have an error  $|\mathcal{B}[\tilde{\mathbf{h}}] - \mathcal{B}_{\bar{K}}[\tilde{\mathbf{h}}]| = \mathcal{O}(\bar{K}^{-1})$  on invariant circles and islands. In contrast, chaotic trajectories are conjectured to have a convergence rate of  $\mathcal{O}(\bar{K}^{-1/2})$ <sup>10</sup>, the same convergence as expected from a central limit theorem.

If we compose the observable with an invariant circle, we can define the observable using coordinates on the torus. That is, if we let  $\mathbf{h} = \tilde{\mathbf{h}} \circ \mathbf{z}$ , we define the finite-time Birkhoff average of the function  $\mathbf{h}$  at a point  $\boldsymbol{\theta} \in \mathbb{T}^d$  as

$$\begin{aligned} \mathcal{B}_{\bar{K}}[\mathbf{h}](\boldsymbol{\theta}) &= \frac{1}{\bar{K}} \sum_{k=0}^{\bar{K}-1} (\mathbf{h} \circ R_{\omega}^k)(\boldsymbol{\theta}), \\ &= \mathcal{B}_{\bar{K}}[\tilde{\mathbf{h}}](\mathbf{z}(\boldsymbol{\theta})). \end{aligned}$$

Assuming  $\mathbf{h}$  is continuous and  $\omega$  is irrational and rationally independent (i.e.  $R_{\omega}$  is ergodic on  $\mathbb{T}^d$ ), the limit of these averages is independent of the initial point  $\boldsymbol{\theta}$ . The average is equal to a spatial average

$$\mathcal{B}[\mathbf{h}] = \lim_{\bar{K} \rightarrow \infty} \mathcal{B}_{\bar{K}}[\mathbf{h}] = \int_{\mathbb{T}^d} \mathbf{h}(\boldsymbol{\theta}) d\boldsymbol{\theta}, \quad (3)$$

where we note that the Lebesgue measure is the unique invariant measure under  $R_{\omega}$ . In this way, we connect averages of time series to averages over invariant tori.

One application of Birkhoff averages is to find Fourier coefficients of  $\mathbf{h}$ . Consider that

$$\mathbf{h}(\boldsymbol{\theta}) = \sum_{\mathbf{n} \in \mathbb{Z}^d} \mathbf{h}_{\mathbf{n}} e^{2\pi i \mathbf{n} \cdot \boldsymbol{\theta}}.$$

Then, the coefficients  $\mathbf{h}_n$  are determined by

$$\begin{aligned} \mathbf{h}_n &= \int_{\mathbb{T}^d} \mathbf{h}(\theta) e^{-2\pi i \mathbf{n} \cdot \theta} d\theta \\ &= \mathcal{B} [\mathbf{h}(\star) e^{-2\pi i \mathbf{n} \cdot \star}] \\ &= \lim_{\bar{K} \rightarrow \infty} \frac{1}{\bar{K}} \sum_{k=0}^{\bar{K}-1} \mathbf{h}(k\omega) e^{2\pi i k \mathbf{n} \cdot \omega}, \end{aligned} \quad (4)$$

where we use ‘ $\star$ ’ for the argument to be averaged over. A special case is  $\mathbf{n} = 0$ , where the constant Fourier term aligns with the unweighted Birkhoff average  $\mathbf{h}_0 = \mathcal{B}[\mathbf{h}]$ . If  $\tilde{\mathbf{h}}$  is the identity and  $X$  is a Euclidean space, the coefficients  $\mathbf{h}_n = \mathbf{z}_n$  provide a Fourier parameterization of an invariant torus

$$\mathbf{z}(\theta) = \sum_{n \in \mathbb{Z}} \mathbf{z}_n e^{2\pi i \mathbf{n} \cdot \theta} = \sum_{n \in \mathbb{Z}} \mathcal{B} [\mathbf{z}(\star) e^{-2\pi i \mathbf{n} \cdot \omega \star}] e^{2\pi i \mathbf{n} \cdot \theta}.$$

To apply the above process for finding coefficients, one must first obtain the rotation vector. For invariant circles, another application of the ergodic average is to obtain the rotation number  $\omega$ . For instance, if one has access to a symplectic map  $\mathbf{F} : \mathbb{T} \times \mathbb{R} \rightarrow \mathbb{T} \times \mathbb{R}$  of the form

$$x_{t+1} = x_t + y_{t+1} \pmod{1}, \quad y_{t+1} = F_y(x_t, y_t), \quad (5)$$

the average of  $y$  is the rotation number of irreducible circles (i.e., those that wind around the torus). Additionally, if one has access to a point  $\mathbf{x}_0$  inside of an invariant circle  $z \subset \mathbb{R}^2$ , one can often find the rotation number by averaging the winding about that point.

However, the winding process to find rotation numbers has some potential difficulties. The most immediate difficulty is with orbits that have complicated shapes (such as crescent- or banana-like orbits), where the average position lies outside of the circle. Even when a point is known to be within the circle, it must be star-shaped for winding to be successful. A more difficult situation is when  $\mathbf{z}$  is not injective, as can occur when considering delay embeddings.

With some modifications, we can modify the above ideas to analyze islands. Using the fact that islands consist of  $p$  periodic circles, one can similarly define Fourier series for the observable on each island  $\mathbf{h}^{(j)} = \tilde{\mathbf{h}} \circ \mathbf{z}^{(j)}$ . Then, we see that the subsequences associated with each island in the chain can be written as

$$\mathbf{a}_{j+kp} = \sum_{n \in \mathbb{Z}^d} \mathbf{h}_n^{(j)} e^{2\pi i k \mathbf{n} \cdot \omega}. \quad (6)$$

For each  $n$ , the coefficients  $\mathbf{h}_n^{(j)}$  can be written as a finite discrete Fourier series

$$\mathbf{h}_n^{(j)} = \sum_{m=0}^{p-1} \mathbf{h}_{mn} e^{2\pi i m j / p}.$$

Combining this with (6), we find that the full sequence can be written as

$$\mathbf{a}_t = \sum_{n \in \mathbb{Z}^d} \sum_{m=0}^{p-1} \mathbf{h}_{mn} e^{2\pi i t (\mathbf{n} \cdot \omega + m) / p} \quad (7)$$

So, signals associated with islands have two frequency components: the rational frequencies  $m/p$  associated with jumping between islands and the irrational frequencies  $\mathbf{n} \cdot \omega / p$  associated with the rotation number. A consequence is that for the map (5), averaging  $y_t$  will typically return a rational number  $m/p$  associated with the number of islands, rather than  $\omega$ . To find  $\omega$  for islands, a two-step process would then be needed, where first the denominator of the average of  $y_t$  is to identify islands (as performed in Sander and Meiss<sup>10</sup>), and then use a second average of  $\mathbf{F}^p$  is used to determine rotation. In the case of higher-dimensional tori and islands, there are multiple irrational frequencies, making it increasingly difficult to solve for frequencies using winding.

Another difficulty is that the  $\mathcal{O}(\bar{K}^{-1})$  convergence rate of the Birkhoff averages on circles and islands is too slow for most applications. For smooth enough maps and invariant circles, this can be improved via the *weighted Birkhoff average*

$$\mathcal{WB}_{\bar{K}}[\tilde{\mathbf{h}}](\mathbf{x}) = \sum_{k=0}^{\bar{K}-1} w_{k, \bar{K}} (\tilde{\mathbf{h}} \circ \mathbf{F}^k)(\mathbf{x}), \quad (8)$$

where the coefficients  $w_{k, \bar{K}}$  are sampled from a positive function  $w \in C^\infty$  compactly supported on  $[0, 1]$  as

$$w_{k, \bar{K}} = \left( \sum_{j=0}^{\bar{K}-1} w \left( \frac{j+1}{\bar{K}+1} \right) \right)^{-1} w \left( \frac{k+1}{\bar{K}+1} \right).$$

To state the convergence rate theorem of the weighted Birkhoff average, we first summarize our assumptions:

**Hypotheses II.1.** *We assume that the three hypotheses hold:*

- H1. (Smooth bump function) *Let  $w \in C^\infty(\mathbb{R})$  have compact support on  $[0, 1]$  and  $w(x) > 0$  for all  $x \in (0, 1)$ .*
- H2. (Diophantine) *Let  $\omega$  be a rotation vector satisfying a  $(c, \nu)$  Diophantine condition.*
- H3. (Smooth system) *Let  $\mathbf{F} : X \rightarrow X$  be a map where  $X$  is a  $C^M$  manifold and  $\mathbf{F} \in C^M$ . Additionally, let  $X_0 \subseteq X$  be a set where  $\mathbf{F}$  is conjugate to island dynamics with  $\mathbf{z}^{(j)} \in C^M$  (see (1)), rotation vector  $\omega$ , invariant measure  $\mu$  and period  $p \geq 1$ . Finally, let  $\tilde{\mathbf{h}} : X \rightarrow \mathbb{R}^D$  be an observable in  $C^M$ .*

Additionally, we note that the simpler case of an invariant torus can be considered by setting  $p = 1$  in the above setting. The following theorem gives the convergence rate of the weighted Birkhoff average:

**Theorem II.2** (Das and Yorke<sup>9</sup> Thm 3.1). *Let  $m > 1$  be an integer. Under Hypotheses II.1, there is a constant  $C_m$  depending on  $w$ ,  $\tilde{\mathbf{h}}$ ,  $m$ ,  $M$ ,  $\nu$ , and  $p$  but independent of  $\mathbf{x} \in X_0$  such that*

$$\left| \mathcal{WB}_{\bar{K}}[\tilde{\mathbf{h}}](\mathbf{x}) - \int_{X_0} \tilde{\mathbf{h}} d\mu \right| \leq C_m \bar{K}^{-m},$$

provided the ‘smoothness’  $M$  satisfies

$$M > d + m\nu.$$

We have modified Thm. II.2 for this work by assuming  $w \in C^\infty$  and extending the theorem to the case of islands. The proof that the theorem works for islands is a simple extension of the invariant torus case, whereby the frequencies in (7) are used instead of those in the original proof. A less smooth function  $w$  could be considered, but we have not found this to be useful in practice.

When  $M = \infty$  for Thm. II.2, one can obtain constants  $C_m$  for any  $m \in \mathbb{N}$ . That is, the weighted Birkhoff average converges as  $\mathcal{O}(\bar{K}^{-m})$  for all  $m$ . This is useful if one wants to quickly compute Birkhoff averages. In contrast, this result does not hold for chaotic trajectories, so the convergence rate of the weighted Birkhoff average can be used to classify trajectories<sup>10</sup>.

One way of understanding the weighted Birkhoff average is as a filter on the sequence  $\mathbf{a}_t = \tilde{\mathbf{h}}(\mathbf{F}^t(\mathbf{x}))$ . When  $\mathbf{x}$  is on an invariant torus, the sequence of observables has the form

$$\mathbf{a}_t = \mathbf{h}(\omega t) = \sum_{n \in \mathbb{Z}} \mathbf{h}_n \lambda_n^t, \quad \lambda_n = e^{2\pi i n \cdot \omega}. \quad (9)$$

So, the sequence  $\mathbf{a}_t$  is built out of equispaced samples of the trigonometric functions  $\theta \mapsto e^{2\pi i \theta \mathbf{n} \cdot \omega}$ . When we apply the weights  $w_{k, \bar{K}}$  to the sequence and flip the order of summation, we find that

$$\mathcal{WB}_{\bar{K}}[\mathbf{h}] = \sum_{n \in \mathbb{Z}} \mathbf{h}_n q(\lambda_n), \quad q(\lambda) = \sum_{k=0}^{\bar{K}-1} w_{k, \bar{K}} \lambda^k.$$

In this way, the weights act as a filter, where  $w_{k, \bar{K}}$  will preferentially remove frequencies from the signal where the polynomial  $q$  is small. An effective filter will return the mean  $\mathbf{h}_0$ , i.e.  $q(1) = 1$ . The rest of the frequencies contribute to the error, so ideally  $|q(\lambda_n)| \ll 1$  for  $\mathbf{n} \neq 0$ . This is particularly important for the small  $|\mathbf{n}|$  modes that dominate the signal.

In Fig. 2, we see an example of this for an example observable  $h(\theta) = e^{\cos(2\pi\theta)}$ . We sample this signal at equispaced points  $h(\omega t)$  where  $\omega = (\sqrt{5} - 1)/2$  is the golden ratio. The top plot shows this signal and the samples. The middle plot shows a windowed discrete Fourier transform of the length 10000 signal. The peaks of the power of the signal appear regularly at multiples of the frequency  $\omega$ , as is expected from the Fourier form of (9). Note that while the peaks occur at the locations  $n\omega$ , these peaks are not ordered as  $n\omega$  wraps around

the torus. This is a consequence of the sampling being below the Nyquist sampling rate, a fact we do not have any practical control over. In the bottom plot, we plot  $|q(e^{2\pi i \omega \theta})|$  for three potential filters with  $\bar{K} = 11$ :

- The ‘all-ones’ filter  $w_{k, \bar{K}} = 1/\bar{K}$ , used for the regular Birkhoff average.
- The weighted Birkhoff filter  $w_{k, \bar{K}}$  sampled from the window function

$$w = e^{-(t(1-t))^{-1}}.$$

- A ‘tuned’ filter that perfectly eliminates the first  $\lfloor \bar{K}/2 \rfloor$  frequencies, found from the coefficients of the polynomial

$$q_{\text{tuned}} = \prod_{k=1}^{\lfloor \bar{K}/2 \rfloor} \frac{(z - e^{2\pi i \omega k})(z - e^{-2\pi i \omega k})}{(1 - e^{2\pi i \omega k})(1 - e^{-2\pi i \omega k})}.$$

In each case, we can judge the absolute error of the filter’s average by comparing it to the exact average (found via a weighted Birkhoff average with  $\bar{K} = 10^4$  to be 1.266066). From Fig. 2 (bottom), we see the all-ones filter polynomial has relatively large value at the peaks of the spectrum, resulting in the worst error of  $7.11 \times 10^{-2}$ . The weighted Birkhoff filter is much smaller at the spectral peaks with the most mass, resulting in a more accurate average with error  $7.38 \times 10^{-3}$ . The tuned filter does the best by two orders of magnitude, with an error of  $2.72 \times 10^{-5}$ .

We note that while the tuned filter is small at the frequencies that dominate the signal, it is large in between. So, while the tuned filter worked well for this example, it would not work well if applied to a signal with a different value of  $\omega$ , as the polynomial roots are specifically related to multiples of the rotation number.

### III. THE BIRKHOFF REDUCED RESIDUAL EXTRAPOLATION METHOD (BIRKHOFF RRE)

At the end of the previous section, we observed an important property: a filter that is tuned to specific frequencies in a signal can be significantly more effective than an arbitrary bump function. In this section, we will introduce a method to learn such a filter from a trajectory.

In Sec. III A, we introduce a continuous problem for finding such an optimal filter on an ergodic component. Then, we discretize this problem with a Birkhoff average in Sec. III B, which results in a variation of the reduced rank extrapolation (RRE) method<sup>23</sup>. Finally, in Sec. III C we explain how we process the obtained filter to find the rotation number of invariant circles.

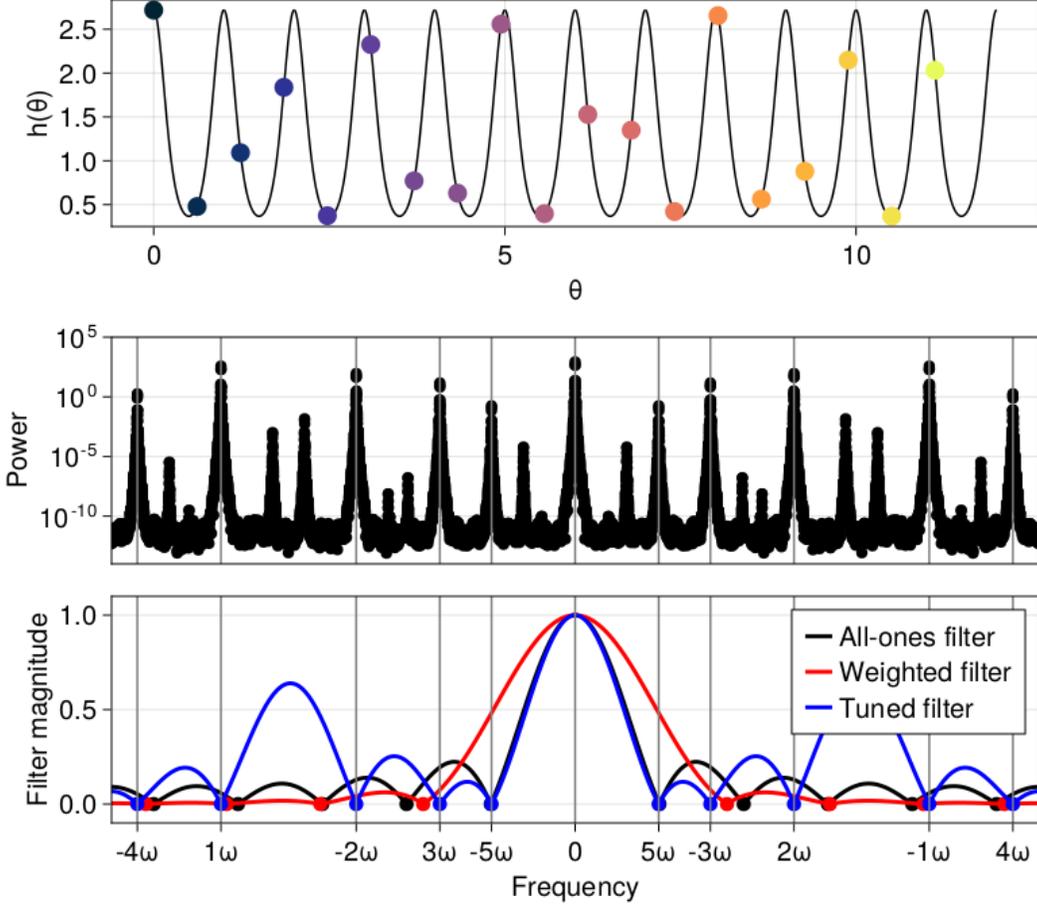


FIG. 2. (top) A test signal  $\mathbf{h}(\theta) = e^{\cos 2\pi\theta}$ , sampled at the points  $\omega t$  where  $\omega = (\sqrt{5} - 1)/2$ . (middle) A discrete Fourier transform of the signal, showing peaks near the expected frequencies. (bottom) Three candidate  $\bar{K} = 11$  filters: the all-ones filter, a weighted Birkhoff filter, and an tuned filter to the first five roots of the frequency. We see that the all-ones filter is largest where there is a large amount of power if the signal, the weighted filter is small far from the zero frequency, and the tuned filter is zero exactly at the relevant frequencies. The absolute errors of the Birkhoff average for each filter are: (all-ones)  $7.11 \times 10^{-2}$ , (weighted Birkhoff)  $7.38 \times 10^{-3}$ , and (tuned)  $2.72 \times 10^{-5}$ .

### A. The Least-Squares Problem

We begin by defining a function for the action of a filter

$$\mathcal{F}_{\bar{K}}[\tilde{\mathbf{h}}](\mathbf{x}_0) = \sum_{k=0}^{\bar{K}-1} c_k (\tilde{\mathbf{h}} \circ \mathbf{F}^k)(\mathbf{x}_0),$$

and we call the associated filter polynomial  $q_{\bar{K}}$ . Note that this is equivalent to a regular Birkhoff average if  $c_k = 1/\bar{K}$  and a weighted Birkhoff average if  $c_k = w_{k,\bar{K}}$ . Our goal is to find coefficients  $c_k$  so that  $\|\mathcal{F}_{\bar{K}}[\tilde{\mathbf{h}}](\star) - \mathbf{h}_0\|_{L^2}$  is small, where we are taking a norm over an ergodic region  $X_0 \subseteq X$  of the form

$$\|f\|_{L^2}^2 = \int_{X_0} |f(\mathbf{x})|^2 d\mu.$$

However, there is a problem: we do not know  $\mathbf{h}_0$  *a priori*, so we cannot directly minimize  $\|\mathcal{F}_{\bar{K}}[\tilde{\mathbf{h}}](\star) - \mathbf{h}_0\|_{L^2}$ .

So, we instead consider filtering the function

$$\tilde{\mathbf{g}}(\mathbf{x}_0) = (\tilde{\mathbf{h}} \circ \mathbf{F})(\mathbf{x}_0) - \tilde{\mathbf{h}}(\mathbf{x}_0).$$

This is a convenient choice because  $\tilde{\mathbf{g}}$  has zero mean and it is easy to calculate from a trajectory of  $\tilde{\mathbf{h}}$ . Given  $\tilde{\mathbf{g}}$ , the new goal is to minimize  $\|\mathcal{F}\tilde{\mathbf{g}}(\cdot)\|_{L^2}$ , under the constraint that  $\sum_k c_k = 1$ . The constraint ensures the associated polynomial satisfies  $q_{\bar{K}}(1) = 1$ , a necessary condition to return the average. Additionally, on an invariant circle or island,  $\tilde{\mathbf{g}}$  has the same type of Fourier series representation as  $\tilde{\mathbf{h}}$ . So, if a filter learns the rotation vector of  $\tilde{\mathbf{g}}$ , it will also be the rotation vector associated with  $\tilde{\mathbf{h}}$ .

To discretize the norm, we use a weighted Birkhoff average (8):

$$\|\mathcal{F}_{\bar{K}}\tilde{\mathbf{g}}\|_{L^2}^2 \approx \mathcal{WB}_T \left[ |(\mathcal{F}_{\bar{K}}\tilde{\mathbf{g}} \circ \mathbf{F}^t)(\star)|^2 \right] (\mathbf{x}_0). \quad (10)$$

Assuming the conditions of Thm. II.2 are met, we can

bound the error of the above approximation by

$$\left| \|\mathcal{F}_{\bar{K}}\tilde{\mathbf{g}}\|_{L^2(\mathbb{T})}^2 - \mathcal{WB}_T \left[ |(\mathcal{F}_{\bar{K}}\tilde{\mathbf{g}} \circ R_\omega^t)(\star)|^2 \right] (\boldsymbol{\theta}_0) \right| < CT^{-m}$$

for some  $C > 0$  and integer  $m$ .

In summary, the sum on the right hand side of (10) can be seen as measure of how good a given filter  $\mathbf{c}$  is on an invariant set. Additionally, one can obtain this approximation using only a single trajectory of the dynamical system, rather than having *a-priori* information. If we had used an unweighted Birkhoff average with no more assumptions, this returns the standard RRE algorithm. However, the weighted Birkhoff average acts as a convenient weighting for RRE by connecting it efficiently to a continuous problem.

## B. Least Squares Solution

Now that we have an energy to minimize, we discuss the numerical details. We begin by observing a structure of invariant circle and island signals: they come from pure Fourier series. That is, there are no growing or decaying modes, so the frequencies that we hope to learn via the filter all correspond to filter polynomial roots on the unit circle. Such roots do not change under time reversal (i.e. the conjugate pair obeys  $(\lambda_n, \bar{\lambda}_n) \rightarrow (\bar{\lambda}_n, \lambda_n)$  as  $t \rightarrow -t$ ). This property corresponds to a linear constraint on the filter that

$$c_{K+k} - c_{K-k} = 0, \quad (11)$$

where  $0 \leq k < K$ . Filters satisfying (11) are known as *palindromic*. We note that the converse is not true: time reversal symmetry does not imply roots on the unit circle. However, while not strictly necessary, we have found that this constraint dramatically improves the quality of the results. Throughout the rest of this paper, we will use  $\bar{K}$  to refer to a filter of length  $2K + 1$ , whereas we used  $\bar{K}$  and  $T$  to represent filters of length  $\bar{K}$  and  $T$  respectively. We will find that  $K$  is the number of unknowns of the final least-squares problem.

For the algorithm in this section, we assume that the user has access to an initial point  $\mathbf{x}_0 \in X$ , a symplectic map  $\mathbf{F}$ , and an observable function  $\tilde{\mathbf{h}}$ . The algorithm begins by sampling a trajectory of length  $T + 2K + 1$  starting at  $\mathbf{x}_0$  by repeated application of  $\mathbf{F}$ . For many test maps, such as the standard map, this step is very quick. However, in applications where evaluating  $\mathbf{F}$  involves simulating a dynamical system up to a Poincaré section, this step could potentially dominate the cost of the algorithm.

From here, we compute the difference sequence  $\mathbf{u}_t = \mathbf{a}_{t+1} - \mathbf{a}_t$  for  $0 \leq t < T$ . This step amounts to computing the difference function  $\tilde{\mathbf{g}} \circ \mathbf{F}^t$  from the previous section. Using this notation, the weighted Birkhoff average in (10) can be expressed as the product

$$\mathcal{WB}_T \left[ |(\mathcal{F}_{\bar{K}}\tilde{\mathbf{g}} \circ \mathbf{F}^t)(\star)|^2 \right] (\mathbf{x}_0) = \mathbf{c}^T U^T W_T U \mathbf{c},$$

where  $U \in \mathbb{R}^{TD \times (2K+1)}$  is the block-Hankel matrix

$$U = \begin{pmatrix} \mathbf{u}_0 & \mathbf{u}_1 & \dots & \mathbf{u}_{2K} \\ \mathbf{u}_1 & \mathbf{u}_2 & \dots & \mathbf{u}_{2K+1} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{u}_{T-1} & \mathbf{u}_T & \dots & \mathbf{u}_{T-1+2K} \end{pmatrix},$$

and  $W_T \in \mathbb{R}^{TD \times TD}$  is a diagonal matrix with the weighted Birkhoff weights

$$W_T = \begin{pmatrix} w_{0,T} I_D & & & \\ & \ddots & & \\ & & & w_{T-1,T} I_D \end{pmatrix},$$

where  $I_D$  is the identity matrix in  $\mathbb{R}^{D \times D}$ . We note that the matrix-vector product  $U\mathbf{c}$  can be interpreted as the filter being applied to sliding windows of  $\mathbf{u}_t$ , i.e.  $(U\mathbf{c})_t = \mathcal{F}_K[\tilde{\mathbf{g}} \circ \mathbf{F}^t](\mathbf{x}_0)$ .

For the full least squares problem, we require two more components. First, a filter must return the correct mean, which corresponds to the constraint that  $\mathbf{c} \cdot \mathbf{1} = q_K(1) = 1$ . Second, we allow for a small amount of regularization to remove the possibility of low rank systems that arise from periodic orbits and very smooth circles. For the regularization, we choose weights so that the solution of the regularized problem is a weighted Birkhoff average. This is performed via the inverse of another weighted Birkhoff matrix  $W_K \in \mathbb{R}^{(2K+1) \times (2K+1)}$  with entries

$$(W_K)_{kk} = \tilde{w}_{k,K} = \frac{w_{k,2K+1} + w_{K-k,2K+1}}{2}.$$

Note that  $W_K$  has been symmetrized to respect the palindromic symmetry, and can equivalently be thought of as sampling the symmetric bump function  $w(1/2 - x) + w(1/2 + x)$ .

In total, we define the following least-squares problem for finding  $\mathbf{c}$ :

$$R^2 = \min_{\mathbf{c} \in \mathbb{R}^{2K+1}} \left\| \begin{pmatrix} W_T^{1/2} U \\ \epsilon^{1/2} W_K^{-1/2} \end{pmatrix} \mathbf{c} \right\|^2, \quad (12)$$

$$\text{s.t. } \mathbf{1} \cdot \mathbf{c} = 1, \quad c_{K+k} - c_{K-k} = 0 \text{ for } 1 \leq k \leq K,$$

where  $\epsilon \geq 0$  is a regularization parameter. The above problem can be recognized as a weighted and time-reversal constrained version of RRE.

We enforce the palindromic constraint by projecting by a matrix  $P \in \mathbb{R}^{K+1 \times 2K+1}$ , where  $\tilde{\mathbf{c}} = P\mathbf{c}$  with

$$(P\mathbf{c})_k = \begin{cases} c_K, & k = 0, \\ \frac{c_{K+k} + c_{K-k}}{\sqrt{2}}, & 1 \leq k \leq K. \end{cases} \quad (13)$$

Substituting the matrix into (12), we find

$$R^2 = \min_{\tilde{\mathbf{c}} \in \mathbb{R}^{K+1}} \tilde{\mathbf{c}}^T \tilde{A} \tilde{\mathbf{c}} + \epsilon \tilde{\mathbf{c}}^T \tilde{W}_K^{-1} \tilde{\mathbf{c}}, \quad (14)$$

$$\text{s.t. } P\mathbf{1} \cdot \tilde{\mathbf{c}} = 1,$$

where

$$\tilde{A} = PU^T W_T U P^T, \quad \tilde{W}_K = PW_K P^T.$$

The effect of multiplying  $U$  by  $P^T$  is to “fold  $U$  in half” and to sum the matching columns.

Computationally, we can fully handle the constraints by including  $\mathbf{1} \cdot \mathbf{c} = 1$  in the projection. We do this by changing variables to a vector  $\boldsymbol{\xi} \in \mathbb{R}^K$  (see Sidi<sup>23</sup>, pp 41)

$$c_k = P_\xi \boldsymbol{\xi} + \mathbf{e}_K,$$

where  $\mathbf{e}_K$  is the  $K$ th unit vector and

$$P_\xi = \begin{pmatrix} 1 & -1 & & & & -1 & 1 \\ & \ddots & \ddots & & & \ddots & \ddots \\ & & 1 & -1 & & -1 & 1 \\ & & & 1 & -2 & 1 & \end{pmatrix}.$$

Clearly, we have  $P_\xi \mathbf{1} = 0$  and that  $(P_\xi)_{k,K+\ell} = (P_\xi)_{k,K-\ell}$ , so every constraint is satisfied. We can substitute this into our least-squares problem, giving

$$R^2 = \min_{\boldsymbol{\xi} \in \mathbb{R}^K} \left\| \begin{pmatrix} W_T^{1/2} U \\ \epsilon^{1/2} W_K^{-1/2} \end{pmatrix} P_\xi \boldsymbol{\xi} + \begin{pmatrix} W_T^{1/2} U \mathbf{e}_0 \\ \epsilon^{1/2} W_K^{-1/2} \mathbf{e}_0 \end{pmatrix} \right\|^2. \quad (15)$$

We solve this system via a direct QR-based least squares solve, as this is reliably accurate and typically fast enough for the system sizes considered.

We note that the system could alternatively be solved via an iterative least-squares solver such as LSQR<sup>26</sup>. The matrices  $P_\xi$ ,  $W_T^{1/2}$ , and  $W_K^{-1/2}$  can all be applied in  $\mathcal{O}(T)$  time. Additionally, the matrix  $U$  can be applied  $\mathcal{O}(T \log T)$  time via fast Hankel multiplications using the fast Fourier transform. So, an iterative algorithm would have a worst-case run time  $\mathcal{O}(TK \log T)$  operations, rather than the  $\mathcal{O}(TK^2)$  time to perform the QR-based algorithm. Additionally, in the case that  $K$  is chosen much larger than necessary, the iterative algorithm converges in many fewer steps. However, we have not yet found an appropriate preconditioner for this method: a necessary step for convergences to high tolerances.

In total, an algorithm for finding a filter is given below:

---

#### Algorithm 1 Reduced Rank Extrapolation

---

**Input:** Initial point  $\mathbf{x}_0$ , symplectic map  $\mathbf{F}$ , least squares dimensions  $T$  and  $K$ , regularization  $\epsilon$

- 1: Sample trajectory  $\mathbf{a}_t$  via repeated evaluations of  $\mathbf{F}$
- 2: Compute  $\mathbf{u}_t$  and weights for  $W_T$  and  $W_K$
- 3: Solve (15) for  $\boldsymbol{\xi}$  via direct QR or iterative LSQR solver
- 4: Compute  $\mathbf{c} = P_\xi \boldsymbol{\xi}$  and  $R$  from (15)

**Output:** Best filter  $\mathbf{c}$  and residual  $R$

---

The above algorithm directly inherits convergence rates of the weighted Birkhoff average.

**Theorem III.1.** *Under Hypotheses II.1, the residual of (12) converges independent of how  $T$  depends on  $K$  as*

$$\epsilon \leq R^2 \leq \epsilon + C_m (2K)^{-2m}, \quad (16)$$

where  $m$  and  $C_m$  are identical to Thm. II.2.

*Proof.* Consider the solution  $c_k = \tilde{w}_{k,K}$ . By construction, the filter obeys the constraints. Using Thm. II.2,  $|(U\mathbf{c})_t| \leq C_m (2K+1)^{-m}$ . So,

$$\mathbf{c}^T U^T W_T U \mathbf{c} = \sum_{t=0}^{T-1} w_{t,T} (U\mathbf{c})_t^2 \leq C_m (2K)^{-2m}.$$

Similarly, the regularization term evaluates exactly as  $\mathbf{c}^T W_K^{-1} \mathbf{c} = \epsilon$ . Combining these computations, we have the upper bound.

For the lower bound, we consider  $\min_{\mathbf{c}} \mathbf{c}^T W_K^{-1} \mathbf{c} \leq R$ . Enforcing the  $\mathbf{1} \cdot \mathbf{c} = 1$  constraint with a Lagrange multiplier  $\lambda$  and the palindromic constraint explicitly, the minimizer is found to solve the linear system

$$\tilde{w}_{k,K}^{-1} c_k + \lambda = 0, \quad \sum_k c_k = 1.$$

This is solved via  $c_k = \tilde{w}_{k,K}$  and  $\lambda = -2$ , the same solution as considered for the upper bound. Substituting  $c_k$  back into the least-squares problem gives the lower bound.  $\square$

From (16), we see that the extrapolation method converges at least as fast as the weighted Birkhoff average, particularly when  $\epsilon = 0$ . This means that we can also use the convergence of RRE to distinguish between chaotic and non-chaotic trajectories.

In the case that  $d = 1$  and  $\epsilon = 0$ , the convergence rate can be improved using a continued fraction argument instead of relying on the weighted Birkhoff average:

**Theorem III.2.** *Let  $0 < \eta < 1$ ,  $w : [0, 1] \rightarrow \mathbb{R}$  be a nonzero positive bounded function, and assume Hypothesis II.1 (H3) with  $d = 1$  so that  $\mathbf{h} \in C^M$ . Then, for almost all  $\omega \in \mathbb{T}$ , there exists an  $L$  such that for  $2K+1 > L$  such that for some  $C > 0$*

$$R^2 < C(2K+1)^{2\eta(-M+1)}.$$

Furthermore, if  $\mathbf{h} \in C^\omega$  is a real analytic function on the torus, then for some  $0 < r < 1$  independent of  $\eta$  and  $C > 0$ , the error obeys the inequality

$$R^2 < C r^{(2K+1)^\eta}.$$

*Proof.* See Appendix C.  $\square$

**Remark III.3.** *We note that the proof of this theorem does not rely on any assumptions on  $w$ . That is, so long as  $W_T$  are positive diagonal matrices with  $\text{Tr}(W_T) = 1$ , the bound above will hold. Additionally, the nearly exponential rate obtained for analytic functions is faster than any known guarantee for the weighted Birkhoff average. This is due to the lack of reliance on bump functions, which can be  $C^\infty$  but cannot be analytic.*

Because neither of the above theorems depend on  $T$ , there are not any strong lower bounds on the residual. Practically, it is required that  $Td$  be greater than or equal

to  $K$  for a full rank linear system. This is an important requirement to ensure that chaotic trajectories do not converge. However, to our knowledge, there is currently no theorem here or in the weighted Birkhoff average literature proving that chaotic trajectories cannot be accelerated, despite strong numerical evidence. We will similarly give no proofs for the extrapolation method in chaos.

In order to classify trajectories efficiently, we also consider an adaptive algorithm. For this, we first define the scale-free residual  $R_G$  as

$$R_G = \frac{\sqrt{R^2 - \epsilon}}{G}, \quad G^2 = \sum_{t=0}^{T-1} w_{t,T} |\mathbf{u}_t|^2 = \mathcal{WB}[\tilde{\mathbf{g}}^2](\mathbf{x}_0). \quad (17)$$

Because  $G$  converges to a nonzero number with  $T$  (assuming  $\mathbf{x}_0$  is not a fixed point),  $R_G$  converges at the same rate as  $R$ . For the adaptive algorithm, we increase  $K$  from an initial value  $K_{\min}$  to a maximum value  $K_{\max}$  by an increment of  $\Delta K$ . We assume that  $T$  scales with  $K$  as  $T = \lceil \gamma K / D \rceil$  for some constant  $\gamma \geq 1$ . At each increment, we check whether  $R_G < \delta$ , and finish the algorithm when this conditions is met. The fact that  $R_G$  is scale-free allows for the algorithm to be applied to multiple orbits for the same map with an accuracy to match the largest scale of the system. In total, the adaptive algorithm is below:

---

**Algorithm 2** Adaptive Reduced Rank Extrapolation

---

**Input:** Initial point  $\mathbf{x}_0$ , symplectic map  $\mathbf{F}$ , regularization  $\epsilon$ ,  $T$  proportionality constant  $\gamma$ , convergence cutoff  $\delta$ , initial  $K$  value  $K_{\text{init}}$ , maximum  $K$  value  $K_{\text{max}}$ ,  $K$  increment  $\Delta K$

- 1:  $K \leftarrow K_{\text{init}}, T \leftarrow \lceil \gamma K / D \rceil$
- 2: **while**  $K \leq K_{\text{max}}$  and  $R_G > \delta$  **do**
- 3:   Get  $\mathbf{c}$ ,  $R$  via Algorithm 1, reusing trajectory information
- 4:   Get  $R_G$  via (17)
- 5:    $K \leftarrow K + \Delta K, T \leftarrow \lceil \gamma K / D \rceil$
- 6: **end while**

**Output:** Best filter  $\mathbf{c}$  and residual  $R$

---

With the use of  $QR$  factorization updates, the adaptive algorithm maintains an  $\mathcal{O}(TK^2)$  runtime. However, even without updates, choosing large enough steps  $\Delta K$  mostly avoids the cost of adaptivity even without updates for reasonably small systems.

### C. Processing the Learned Filter

Once we have found a filter  $\mathbf{c}$  for the sequence of observations for a trajectory, we must post-process the results. There are two steps to this process. The first is straightforward: we distinguish chaotic trajectories from non-chaotic trajectories via a tolerance for the residual. The same tolerance as in Algorithm 2 can be used, although often it is convenient to use separate tolerances for adaptivity convergence and chaos classification. If the

residual is above the tolerance, we deem the trajectory “chaotic,” and otherwise it is an invariant circle or an island. We note that choosing the tolerance for a specific problem depends on the desired accuracy, the tolerable amount of computation, and the accuracy of the  $\mathbf{F}$  evaluation. However, any error in the computation of  $\mathbf{c}$  will lead to more difficulty in the following steps, so lower tolerances are always more reliable for adaptive RRE than classification.

In the case that we deem the trajectory non-chaotic, the second step is to process the roots of the filter polynomial. Empirically, we have found the following convergence result to hold for the polynomial roots:

**Conjecture III.4.** *Under Hypotheses II.1, let  $\mathbf{c}_K$  be a sequence of solutions to (14) with  $T = \lceil \gamma K \rceil$  for  $\gamma \geq 1$  and let  $q_K(z)$  be the filter polynomial with coefficients  $\mathbf{c}_K$ . Then, for all  $n$  such that  $|\mathbf{g}_n| \neq 0$ , there exists an  $M_* > 0$  and a sequence  $z_{K,n}$  such that  $q_K(z_{K,n}) = 0$  and*

$$|\lambda_n - z_{K,n}| < C_n K^{-M_*}, \quad (18)$$

where  $M_*$  depends on  $M$  and potentially  $\nu$ . In particular, if  $M = \infty$ ,  $M_*$  is unbounded.

The above conjecture essentially states that the roots of the filter necessarily approach the points  $\lambda_n$  associated with multiples of the rotation number and the island period. The difficulty with in proving the conjecture is that the weighted Birkhoff average performs strongly for minimizing the RRE residual. So, to prove the roots converge, one must show that filter polynomials with roots that converge to  $\lambda_n$  out-perform weighted Birkhoff averages in minimizing the residual.

Numerically, roots of the filter polynomial are found by solving the equation

$$q_K(\lambda) = \sum_{k=0}^{2K} c_k \lambda^k = 0,$$

for every root  $\tilde{\lambda}_j$ . This is performed by reducing the polynomial to another Chebyshev polynomial with half the dimension, then solving for the eigenvalues of the associated colleague matrix<sup>27</sup> (see App. A). The dimension reduction accelerates the root-finding by a significant constant factor. Additionally, due to the majority of eigenvalues being on the unit circle, this eigenvalue problem is well conditioned. The only exception to this is when  $\lambda = -1$  is a root, as happens with  $p = 2$  islands. In this case, the palindromic coefficients force  $-1$  to be a double root, causing a square root of the error of this frequency.

Throughout the rest of this section, we assume that  $d = 1$  and therefore the trajectory is either an invariant circle or a  $d = 1$  island. Using the roots  $\tilde{\lambda}_n$ , we will find the rotation number by finding which roots represent the majority of the signal. This step is heuristic: we assume that the low-frequency oscillations will dominate the Fourier spectrum, so we can use this information to

prune the important frequencies from the less important ones. To do this, we first restrict to values of  $\lambda_n$  that are very close to the unit circle by some tolerance on the order of the square root of machine epsilon (due to the aforementioned sensitivity of  $-1$  as a root). The remaining eigenvalues are sorted by solving the weighted least-squares system for each eigenvalue's prominence in the signal

$$\min_V \left\| W_{2K+T+1}^{1/2} (\tilde{\Phi}V - A) \right\|^2,$$

where  $\tilde{\Phi} \in \mathbb{R}^{2K+T+1 \times 2K+1}$  with  $\tilde{\Phi}_{mn} = \tilde{\lambda}_n^m$  is the matrix of eigenmodes associated with the frequencies  $\tilde{\lambda}_n$ ,  $V \in \mathbb{R}^{2K+1 \times D}$  holds the prominence of each mode in the signal,  $W_{2K+T+1} \in \mathbb{R}^{2K+T+1 \times 2K+T+1}$  is the weighted Birkhoff diagonal matrix with  $(W_{2K+T+1})_{mm} = w_{m,2K+T+1}$ , and

$$A = \begin{pmatrix} \mathbf{a}_0^T \\ \vdots \\ \mathbf{a}_{T+2K}^T \end{pmatrix}.$$

The eigenvalues are then sorted by the value of norm of the rows of  $V$ .

Using the top few eigenvalues (arbitrarily chosen to be 10), we next determine if the sequence corresponds to an island. We do this by noticing a key difference between a  $p \geq 2$  island and a Diophantine invariant circle: the island spectrum (7) has rational frequencies whereas the invariant circle (9) does not. So, if the sequence is an island, we expect a prominent rational frequency  $\tilde{\omega}_n = (2\pi i)^{-1} \arg \tilde{\lambda}_n$  in the sorted eigenvalues. For this, we use a Farey-sequence method from Sander and Meiss<sup>10</sup> (Algorithm 2) to determine whether a given root is rational to a specified tolerance  $\epsilon_{\text{rat}}$  and bounded denominator  $p \leq p_{\text{max}}$ . If one of the roots is rational by this algorithm, we consider it to correspond to an island chain with the largest found period  $p$ . From here, we can rerun Algorithm 1 on the ‘‘stacked’’ signal with dimension  $\hat{D} = pD$

$$\hat{\mathbf{a}}_t = \begin{pmatrix} \mathbf{a}_{pt} \\ \mathbf{a}_{pt+1} \\ \vdots \\ \mathbf{a}_{pt+p-1} \end{pmatrix},$$

with a filter length  $\hat{K} = \lfloor K/D \rfloor$ . Using this signal is equivalent to working simultaneously with the  $p$  invariant circles of  $\mathbf{F}^p$  with the same rotation number. The resulting filter will be that of an invariant circle and the above steps can be repeated, skipping the check for rational roots.

Now, consider the case that no rational roots of the filter are found. In order to determine the rotation number, we simply choose the frequency that appears first in the sorted eigenvalues, i.e.  $\omega = (2\pi i)^{-1} \arg(\tilde{\lambda}_0)$  where  $\tilde{\lambda}_0$  is the first eigenvalue after sorting. We note that this

is not always true: invariant circles can be found where the largest Fourier coefficient does not correspond to the rotation number. However, we have not found this to be a practical issue for any of our examples.

Once the rotational frequency  $\omega$  is determined, we can obtain the coefficients of the invariant circle associated with the signal by solving another linear system. This time, because typically  $\omega$  is resolved to very high accuracy, we solve the linear system

$$\min_V \left\| W_{2K+T+1}^{1/2} (\Phi V - A) \right\|^2, \quad (19)$$

where  $\Phi \in \mathbb{R}^{2K+T+1 \times L}$  with  $\Phi_{mn} = \lambda_n^m$ ,  $\lambda_n = e^{2\pi i \omega(n-L)}$  are the new signal frequencies, and  $L$  is the number of coefficients wanted, chosen to control the condition number of the least-squares problem (see Sec. B). The approximated invariant circles are then given as

$$\mathbf{z}^{(j)}(\theta) = \sum_{\ell=-L}^L V_{\ell, D(j-1):Dj-1} e^{2\pi i \ell \theta},$$

where we are using ‘‘Matlab notation’’ for slicing the rows of  $V$  into each island component.

Once an invariant circle or island chain is found, the fit can be validated via a residual from the parameterization. In particular, we can evaluate the residual

$$R_p^2 = \frac{1}{pJ} \sum_{j=0}^{J-1} \left( \left| \mathbf{z}^{(1)}(jh + \omega) - \mathbf{F}(\mathbf{z}^{(p)}(jh)) \right|^2 + \sum_{j=1}^{p-1} \left| \mathbf{z}^{(j+1)}(hj) - \mathbf{F}(\mathbf{z}^{(j)}(hj)) \right|^2 \right), \quad (20)$$

$$\approx \frac{1}{p} \left( \int_{\mathbb{T}} \left| \mathbf{z}^{(1)}(\theta + \omega) - \mathbf{F}(\mathbf{z}^{(p)}(\theta)) \right|^2 d\theta \right) \quad (21)$$

$$+ \sum_{j=1}^{p-1} \int_{\mathbb{T}} \left| \mathbf{z}^{(j+1)}(\theta) - \mathbf{F}(\mathbf{z}^{(j)}(\theta)) \right|^2 d\theta \right), \quad (22)$$

where  $h = 1/J$  for  $J \in \mathbb{N}$ . The residual  $R_p$  is essentially an  $L^2$  measurement of the conjugacy (1). If this residual is small, then the conjugacy is likely correct. Furthermore, a small  $R_p$  likely puts the island within the basin of convergence for the parameterization method<sup>4</sup>, which could be used to refine the estimate of the island. Thus, Birkhoff RRE could be seen as a method of finding an initial guess for higher-accuracy methods.

#### IV. EXAMPLES

In this section, we give two methods of the Birkhoff RRE method applied to symplectic maps. In Sec. IV A,

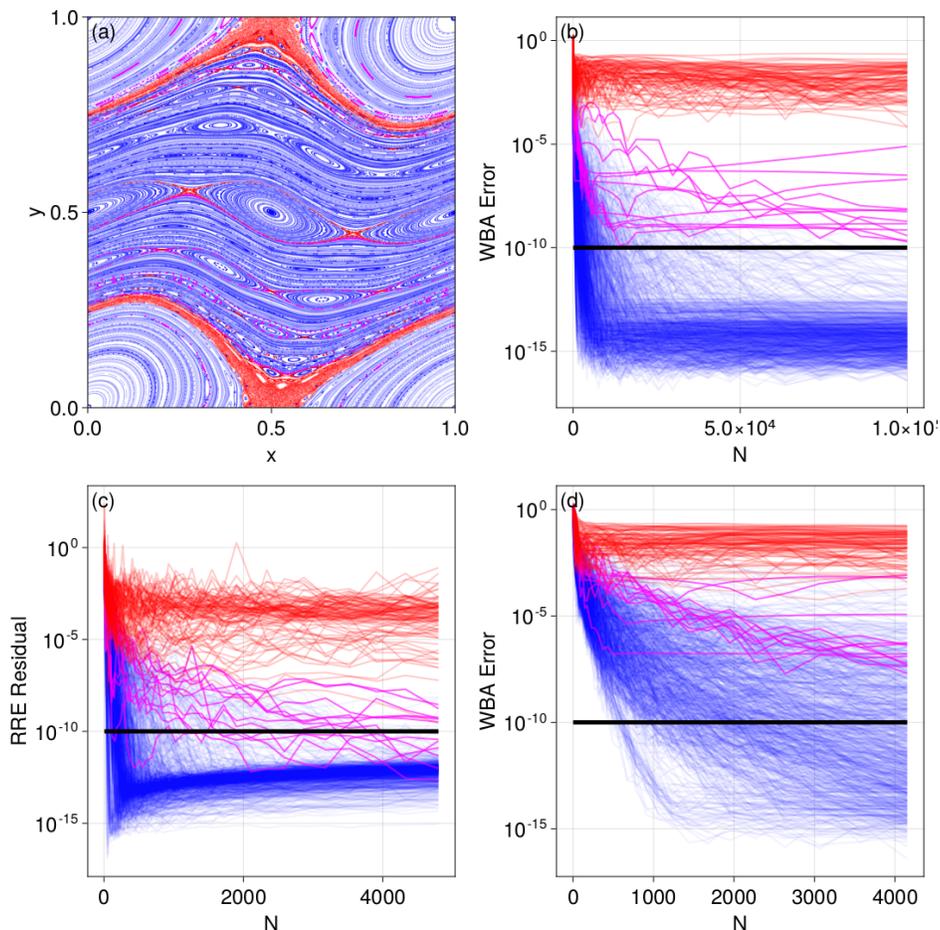


FIG. 3. (a) A Poincaré plot of the standard map with  $k = 0.7$ , colored by trajectory classification of integrable, chaotic, or indeterminate. (b) Convergence of the weighted Birkhoff doubling error with respect to trajectory length. Trajectories are colored by the ending point on this plot, with  $R > 10^{-5}$  chaotic,  $R < 10^{-11}$  integrable, and the rest indeterminate. (c) Convergence of RRE residual with respect to trajectory length. (d) The same as data as (b), on the same domain as (c). From (c) and (d), we see that the RRE residual appears to converge much more rapidly than the weighted Birkhoff average.

we apply Birkhoff RRE to the Chirikov standard map in order to investigate its convergence properties. We show that the convergence of the RRE residual is significantly more efficient in the number of samples than the weighted Birkhoff average (WBA). As a consequence, RRE can classify trajectories with many fewer symplectic map iterations. We also show that the roots of the filter polynomial converge like Conjecture III.4.

In the Sec. IV B, we explore the performance of Birkhoff RRE on magnetic field-line dynamics for a stellarator, a type of plasma confinement device. We show how the method classifies chaotic and integrable trajectories, as well as finding Fourier representations of circles and islands. The rotation number and residuals are also reported. The symplectic map is obtained by numerically integrating magnetic field lines from an interpolated field, so this example additionally shows the performance on a map with symplecticity breaking error.

Code for performing both Birkhoff RRE and the

weighted Birkhoff average herein can be found in the `SymplecticMapTools.jl`<sup>28</sup> Julia package.

### A. Standard Map Convergence

For the first experiment, we are interested in comparing the convergence and classification of Birkhoff RRE vs WBA. For RRE, we measure convergence via the square root of the least-squares residual  $R$  defined in (14). For WBA, we use the method of Sander and Meiss<sup>10</sup>, and compare the values of two consecutive averages. That is, given a symplectic map  $\mathbf{F}$  and an observable  $\tilde{\mathbf{h}}$ , the WBA residual is

$$\begin{aligned}
 R_{\text{WBA}} &= \left\| \mathcal{WB}_T[\tilde{\mathbf{h}}](\mathbf{x}_0) - \mathcal{WB}_T[\tilde{\mathbf{h}}](\mathbf{F}^T(\mathbf{x}_0)) \right\|, \\
 &= \left\| \sum_{t=0}^{T-1} w_{t,T} \mathbf{a}_t - \sum_{t=0}^{T-1} w_{t,T} \mathbf{a}_{t+T} \right\|.
 \end{aligned}$$

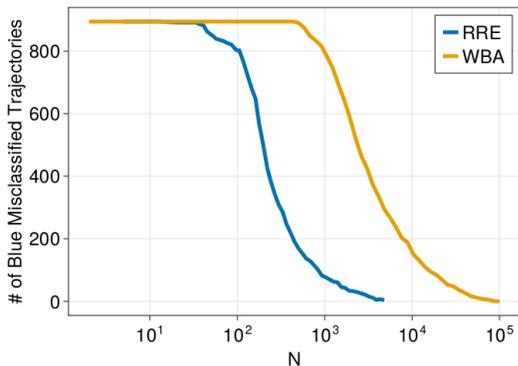


FIG. 4. Number of integrable trajectories of Fig. 3 misclassified, with the tolerance of both WBA and RRE set to  $10^{-11}$ . We see that RRE converges to a low misclassification rate much more efficiently than WBA in the number of map iterations.

Both residuals have the same theoretical rate guarantees for  $C^\infty$  functions, so it is necessary to numerically check that RRE converges significantly faster.

For a fair comparison, we hold the total number of map evaluations  $N$  constant between the two methods. In particular, we compare  $N = (2 + \gamma)K + 1$  for RRE against  $N = 2T$  for WBA, where we choose the ‘‘rectangularity’’ constant to be  $\gamma = 2$ . The example dynamical system we compute with is the Chirikov standard map (2) with  $k_{\text{sm}} = 0.7$  (see also Fig. 1). In order to handle the map domain  $\mathbb{T} \times \mathbb{R}$ , we use the smooth observable  $\tilde{h} : \mathbb{T} \times \mathbb{R} \rightarrow \mathbb{R}^2$

$$\tilde{h}(x, y) = (y + 0.5) \begin{pmatrix} \cos(2\pi x) \\ \sin(2\pi x) \end{pmatrix}.$$

A Poincaré plot of the standard map is found in Fig. 3 (a), where 1000 trajectories are plotted. The trajectories are classified by a long-time  $N = 10^5$  weighted Birkhoff average, with the three categories of nested invariant circles and islands (blue,  $R_{\text{WBA}} < 10^{-11}$ ), chaos (red,  $R_{\text{WBA}} > 10^{-5}$ ), and indeterminate (purple,  $10^{-5} \leq R_{\text{WBA}} \leq 10^{-11}$ ). The convergence of the WBA residual used to classify these trajectories as a function of  $N$  is found in Fig. 3 (b). Visually, the ‘indeterminate’ trajectories are typically either at the transition of trajectory types (e.g. circles to islands or chaos) or have a near small-denominator rotation number. We note that difficult classification on transitional domains is the typical behavior of any numerical method depending on a continuous quantity. Additionally, misclassifications are likely to happen for any trajectory classification algorithm, as orbits on one side of a classification boundary can shadow orbits on the other side for arbitrarily long periods of time.

In Fig. 3 (c), we plot the convergence of the RRE residual for  $N \leq 2801$  and  $\epsilon = 0$ . For comparison, we plot the WBA residual on the same domain in Fig. 3 (d). We see that RRE converges significantly quicker than WBA,

with most residuals reaching the machine precision limit before  $N = 1000$ .

In Fig. 4, we compare the number of misclassified integrable trajectories for RRE and WBA. To obtain this number, we subtract the number of blue trajectories below the  $10^{-11}$  tolerance (shown as a black horizontal line in Figs. 3 (b)-(d)) from the total number of blue trajectories for each value of  $N$ . We see that RRE classification happens an order of magnitude faster than that of WBA.

Next, we consider the convergence of the roots of the filter polynomials. For this, we choose three invariant circles of the standard map of varying smoothness, plotted in Fig. 5 (a). The inner circle is near the core of the nested circles, and is well approximated by a small number of Fourier modes. The yellow circle is more complex, and the green circle is a case on the edge of chaos.

For each circle, we perform RRE to high accuracy ( $K = 1500$ ,  $\gamma = 3$ ) and obtain the rotation number  $\omega$  by the process described in Sec. III C. Then, for increasing  $K$ , we let  $q_K$  be the filter polynomial for  $\gamma = 3$  and let  $z_{K,n}$  be the associated roots. From Conjecture III.4, we expect the roots to converge faster than  $\mathcal{O}(K^{-M_*})$  for all  $M_*$ . To measure this, in Fig. 5 (b-d) we plot the error

$$E_{K,m} = \min_n |z_{K,n} - e^{2\pi i m \omega}|$$

for a variety of values of  $m$  for each circle. In all cases, we find that the convergence for lower multiples of the rotation number is faster than that for higher multiples. This likely due to the higher prominence of these Fourier modes in the signal. For both the small and medium circles, the rotation number converges to the correct value for values of  $K$  less than 100, corresponding to less than 500 iterations of the map. For the outer circle, the roots converge significantly more slowly, reaching machine precision at around  $K = 400$ . Additionally, the lines are approximately straight, indicating a nearly exponential convergence of  $\mathcal{O}(e^{-\alpha K})$  for some  $\alpha > 0$ . As the circles become more difficult to approximate, higher multiples of the rotation number are found to converge due to a higher number of modes with nontrivial Fourier coefficients, with over 50 multiples converging to high accuracy for the outer circle.

## B. A Stellarator Example

In this section, we consider our method applied to stellarators, a type of toroidal plasma confinement device. Stellarators use large magnetic fields to confine the charged particles that comprise a plasma within the device. In particular, high energy particles approximately follow magnetic field lines, described by the dynamical system

$$\dot{\mathbf{x}} = \mathbf{B}(\mathbf{x}),$$

where  $\mathbf{B}$  is the magnetic field and  $\mathbf{x} = r\mathbf{e}_r(\phi) + z\mathbf{e}_z$  is the cylindrical position vector with radial position  $r$ , azimuthal angle  $\phi$ , and  $z$  is the vertical position. Magnetic

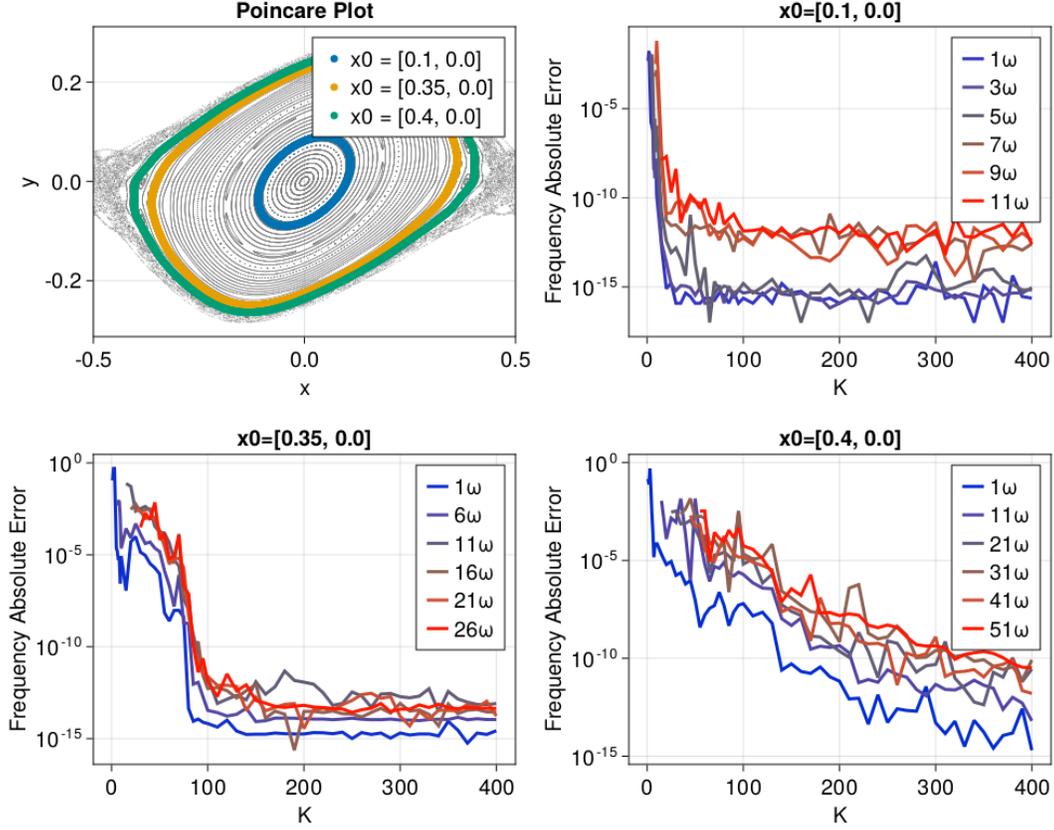


FIG. 5. (a) A Poincaré plot of the standard map, with three invariant circle with initial  $x$  points 0.1, 0.35, and 0.4. (b-d) Convergence plots for the absolute error of the learned frequencies from the polynomial filter vs filter length  $K$ . The absolute error is computed by comparing to the frequency learned from a  $K = 1500$  simulation. In (b), we see very rapid convergence for the smooth central circle. In (c), we see slower, but still rapid convergence of the learned frequencies. In (d), for an invariant circle near the edge of chaos, we see significantly slower convergence. The errors are approximately straight in this plot, indicating nearly exponential convergence of the roots.

field line dynamics is a 1.5D Hamiltonian system, and a symplectic map is obtained by numerically integrating the magnetic ODE over a field period in  $\phi$ . We note that the stellarator map has another difference from the standard map: the domain  $(r, z) \in \mathbb{R}^2$  does not have the rotation number “built in.” In the case of the standard map, an average of the  $y$  coordinate gives the rotation number, whereas the magnetic axis of the stellarator is not necessarily known *a priori*.

To evolve the magnetic field, we use the Julia `Tsit5` integrator on a configuration with many islands. The magnetic field is interpolated from the output of a plasma equilibrium solver using cubic splines. A Poincaré plot showing the 1000 trajectories is found in Fig. 6(a). Due to both the integrator and the field interpolation, there are small non-symplectic errors in the map. These errors affect the clarity of classifying trajectories, as the Birkhoff RRE residuals do not reach the same low levels as the previous example. The errors also make identifying rational numbers and rotation numbers more difficult, and the associated tolerances must also be adjusted. This po-

tentially leads to a increased chance of misclassification, so we perform the following simulation to ensure Birkhoff RRE is robust to errors.

We perform the adaptive Birkhoff RRE classification algorithm 2 on 1000 trajectories using  $\tilde{h} : (r, z) \mapsto (r, z)$ ,  $K_{\min} = 50$ ,  $K_{\max} = 400$ ,  $\Delta K = 50$ ,  $\gamma = 3$ ,  $\epsilon = 0$ , and  $\delta = 10^{-7}$ . With these parameters, a maximum of  $(2+\gamma)K_{\max} + 1 = 2001$  iterations of the map are used per invariant circle, and a minimum of  $(2+\gamma)K_{\min} + 1 = 251$  iterations are used. We note that the cutoff  $\delta$  is chosen low enough to distinguish chaos, but high enough that it is resilient to the noise of evaluating the return map (compare to the value of  $10^{-11}$  for the standard map). After the adaptive RRE algorithm is run, we attempt to fit an invariant circle or island chain to every trajectory, regardless of whether the algorithm converged before the tolerance. Then, we evaluate the validation error 21 for each invariant circle.

We plot the invariant circles and islands parameterizations in Fig. 6(b-c), with the islands colored by the period  $p$ . We only plot those with a validation error below

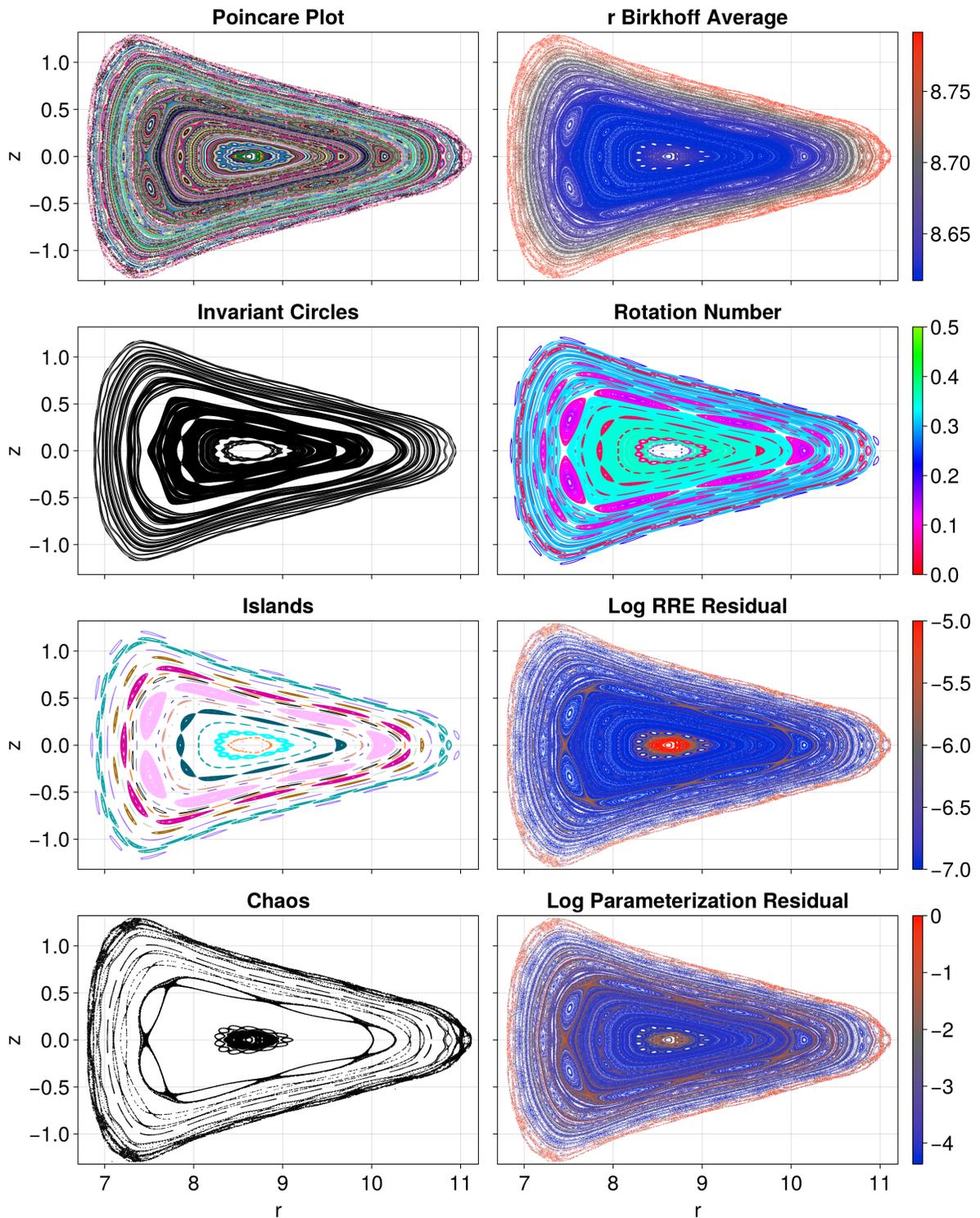


FIG. 6. Eight plots of an optimized stellarator. On the left, we show a Poincaré plot and the trajectories separated into invariant circles, islands, and chaos. On the right, we show the Birkhoff average of the  $r$  coordinate, the rotation number, the log RRE residual, and the log parameterization residual, all computed using the Birkhoff RRE procedure.

the threshold  $R_p < 10^{-2}$ . The associated rotation numbers to the circles and islands are plotted in Fig. 6(f). In Fig. 6(d), we plot the chaotic trajectories, defined as those with the RRE residual satisfying  $R > 5 \times 10^{-7}$ . We note that this value is different than the tolerance used for classification. This is because the standard needed to classify as integrable is separate from the accuracy needed to obtain an effective parameterization.

The circles, islands, and chaos all match expectations from the Poincaré plot. The classification algorithm makes it clear that much of the volume is taken by nested invariant circles, with many islands of various periods intermixed. The classification successfully identifies high-period and low-radius islands that may have been missed without computer processing. Additionally, chaotic trajectories are found at the boundaries of invariant circles and islands, matching the intuition provided by the standard map example. The rotation number plot gives a visual representation of the shear. This allows us to identify the island chains with specific resonant low-denominator frequencies, e.g. with the largest islands occurring at  $\omega = 2/7, 3/11,$  and  $3/10$ .

In Fig. 6(f-g), we plot the Birkhoff RRE residual  $R$  and validation error  $R_p$  on a logarithmic scale. We find the two plots have similar features: both residuals tend to be larger in regions of chaos (the inner core and outer ring) and smaller where there are integrable trajectories. However, we find that the parameterization residual tends to give a sharper indicator of correctness, with chaotic trajectories typically appearing as more red.

## V. CONCLUSIONS

We have shown how Birkhoff RRE, a version of the reduced rank extrapolation algorithm, can be used to accelerate ergodic averages on invariant tori. We find numerically find that this acceleration is significantly more stronger than the weighted Birkhoff average for the same number of map evaluations, at the cost of a more complex algorithm. The acceleration can be used to efficiently classify trajectories in symplectic maps as integrable or chaotic. Beyond this, using the resultant filter from, we can compute the number of islands and the rotation number of integrable trajectories to high accuracy.

Due to Birkhoff RRE post-processing step's reliance on Conjecture III.4, a proof would be a natural future direction from this work. The proof of this convergence would contrast with standard convergence proofs for RRE, which all rely all but finitely many of the frequencies  $\lambda_j$  having modulus less than one. This does not hold for Birkhoff RRE (we have  $|\lambda_j| = 1$  for all  $j$ ), so the proof would likely instead rely on a combination of the Diophantine property and smoothness of  $\mathbf{h}$ , just as both the KAM theorem and the weighted Birkhoff average do.

Beyond this, it would be interesting to apply Birkhoff RRE to higher-dimensional tori. We have proven that the Birkhoff RRE residual converges in higher dimensions, so

it should be able to effectively classify trajectories. Additionally, it is likely that the frequencies converge with similar rates, allowing for high-dimensional Diophantine vectors to be identified without any winding argument. This would be a convenient option for processing high-dimensional symplectic geometry without the need for complicated initial guesses or continuation.

## Appendix A: Finding the Roots of a Palindromic Polynomial

Here, we detail how we solve for the roots of palindromic polynomials. This method is no more or less accurate than solving for the eigenvalues of the standard Frobenius companion matrix (see the book by Trefethen<sup>29</sup> for an introduction). However, this process halves the dimension of the eigenvalue problem, resulting a significantly faster method than the standard companion matrix approach.

The algorithm results from the observation that palindromic polynomials of the form

$$P(z) = \sum_{k=0}^{2K} c_k z^k,$$

can be reduced to degree  $K$  polynomials in  $(z + z^{-1})/2$ . In particular, we have that

$$z^{-K} P(z) = Q\left(\frac{z + z^{-1}}{2}\right), \quad Q(x) = \sum_{k=0}^K b_k T_k(x),$$

where  $T_k$  are the Chebyshev polynomials. If all of the roots of  $P$  are on the unit circle, the change of variables implies that all the roots of the  $Q$  are on the interval  $[-1, 1]$ , so we expect Chebyshev polynomials to be well conditioned for this problem. Typically results of Birkhoff RRE have almost all of the roots on the unit circle for non-chaotic orbits, so they satisfy this assumption. Additionally, using the Chebyshev three-term recurrence

$$\begin{aligned} T_0(x) &= 1, \\ T_1(x) &= x, \\ T_{n+1}(x) &= 2xT_n - T_{n-1}, \end{aligned}$$

one sees that the polynomials in  $z$  coordinates have a convenient form of

$$T_k\left(\frac{z + z^{-1}}{2}\right) = \frac{z^k + z^{-k}}{2}.$$

This results in the convenient relations

$$b_k = \begin{cases} c_K, & k = 0, \\ 2c_{K+k}, & 0 < k \leq K. \end{cases}$$

To find the roots of  $p$ , we use the Chebyshev colleague

matrix

$$C = \begin{pmatrix} 0 & 1 & & & \\ \frac{1}{2} & 0 & \frac{1}{2} & & \\ & \frac{1}{2} & 0 & \ddots & \\ & & \ddots & \ddots & \frac{1}{2} \\ & & & & \frac{1}{2} & 0 \end{pmatrix} - \frac{1}{2b_K} (b_0 \dots b_{K-1}) \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix}.$$

The eigenvalues  $x_n$  of this matrix, found via `eigvals` in Julia, are the roots of the polynomial  $q$ . In order to get roots of  $p$ , we simply solve the equation

$$\frac{z + z^{-1}}{2} = x_n$$

for both values of  $z$ . We note that this equation is sensitive near  $x_n = \pm 1$ , so the Chebyshev method here does not avoid larger errors of frequencies near 0 and  $1/2$ .

### Appendix B: Choosing the Number of Circle Interpolation Modes

Here, we give a method to determine the number of Fourier modes to project a given trajectory onto. We do this by controlling the condition number of the linear system (19). We define the least-squares condition number of the linear system defined by the matrix  $Y = W_T^{1/2} \Phi$  as

$$\kappa = \|Y^\dagger\|_2 = \max_v \frac{\|Y^\dagger \mathbf{v}\|}{\|v\|_2} = \frac{\sigma_1}{\sigma_L},$$

where  $Y^\dagger = (\Phi^T W_T \Phi)^{-1} \Phi^T W_T^{1/2}$  is the Moore-Penrose pseudoinverse and  $\sigma_1$  and  $\sigma_L$  are the largest and smallest singular values of  $Y$ .

To bound the condition number, we first observe that the matrix  $Y^T Y$  has a Toeplitz structure:

$$(Y^T Y)_{mn} = \eta_{n-m} = \sum_{t=0}^T w_{t, T+1} \lambda_1^{(n-m)t}, \quad (\text{B1})$$

where we note that  $\eta_0 = 1$ . In this way, we have a Gershgorin circle bound on the eigenvalues  $\sigma_j^2$  of  $Y^T Y$  of

$$|\sigma_j^2 - 1| \leq 2\gamma_L = 2 \sum_{n=1}^{2L} |\eta_n|,$$

where we used the fact that  $|\eta_n| = |\eta_{-n}|$ . This can be translated to a bound on the condition number of

$$\kappa \leq \frac{\sqrt{1 + 2\gamma_L}}{\sqrt{1 - 2\gamma_L}}. \quad (\text{B2})$$

The expression (B2) allows for an algorithm to bound the condition number of this linear system. We simply choose a maximum allowed radius  $\gamma_{\max}$ , and find the

largest  $L$  such that  $\gamma_L < \gamma_{\max}$ . This requires  $\mathcal{O}(LT)$  operations, which is always cheaper than the linear solve. By default, we have found  $\gamma_{\max} = 0.5$  to give good results. The algorithm is summarized in the following pseudocode:

---

#### Algorithm 3 Parameterization Dimension $L$

---

**Input:** System height  $T$ , Frequency  $\omega$ , Tolerance  $\gamma_{\max}$

```

1:  $L \leftarrow \lceil \frac{T-1}{2} \rceil$ 
2:  $\gamma_0 \leftarrow 0; n \leftarrow 0$ 
3: while  $n \leq T$  do
4:    $n \leftarrow n + 1$ 
5:    $\gamma_n \leftarrow \gamma_{n-1} + \eta_n$  via (B1)
6:   if  $\gamma_n \geq \gamma_{\max}$  then
7:     return  $L \leftarrow \lceil \frac{n-2}{2} \rceil$ 
8:   end if
9: end while
```

**Output:** Parameterization dimension  $L$

---

### Appendix C: Proof of Theorem III.2

In this appendix, we show that for  $d = 1$  and a fixed regularity  $M$ , the RRE residual, and therefore the Birkhoff average, converges at a rate faster than the standard guarantee for the weighted Birkhoff average. This improvement is primarily due to the use of the continued fraction representation of the rotation number rather than the Diophantine property. Using continued fractions, we find progressively better ‘‘ideal polynomials,’’ which exactly filter the leading frequencies of the signal. Because these exactly match the frequencies (rather than the WBA filtering everything evenly), we approach the optimal rate given by the regularity.

The only restriction for the proof is that the continued fraction denominators are sufficiently close to each other, a fact which is given by a standard theorem. Using this, we then show that an ideal polynomial with a degree given by the continued fraction denominator is bounded. This bound gives enough information for the result.

#### 1. The distribution of continued fraction denominators

We begin by quoting a standard theorem on the distribution of continued fraction denominators. The continued fraction representation of a rotation number  $\omega$  is given by the fraction

$$\omega = \frac{1}{c_1 + \frac{1}{c_2 + \frac{1}{\ddots}}}$$

where the coefficients  $c_n \in \mathbb{Z}$  are positive integers. By truncating the series for finitely many  $c_n$ , we obtain the convergents

$$\frac{N_1}{L_1} = \frac{1}{c_1}, \quad \frac{N_2}{L_2} = \frac{1}{c_1 + \frac{1}{c_2}}, \quad \frac{N_3}{L_3} = \frac{1}{c_1 + \frac{1}{c_2 + \frac{1}{c_3}}}, \dots$$

The continued fractions are best approximations of  $\omega$ , satisfying the inequality

$$\left| \omega - \frac{N_n}{L_n} \right| \leq \frac{1}{L_n L_{n+1}} < \frac{1}{L_n^2}. \quad (\text{C1})$$

This is the fundamental result which we will use for the proof of convergence.

Remarkably, one can prove that almost all continued fractions have denominators with similar limiting behavior:

**Theorem C.1.** *For almost all  $\omega \in [0, 1)$ , the denominators of continued fractions  $L_n$  obey the limit*

$$\lim_{n \rightarrow \infty} L_n^{1/n} \rightarrow \gamma,$$

where

$$\gamma = \exp\left(\frac{\pi^2}{12 \log 2}\right).$$

For an introduction to this theorem and more on the measure theory of continued fractions, a good reference is the book by Khinchin<sup>30</sup>.

We then transform this theorem to bounds on continued fraction denominators.

**Corollary C.2.** *Let  $r, s > 0$  and  $\eta < 1$ . For almost all  $\omega \in [0, 1)$ , there is a value  $M$  such that for all  $m > M$ , there exists a continued fraction approximation  $p_n/L_n$  of  $\omega$  such that*

$$rm^\eta < L_n < sm. \quad (\text{C2})$$

*Proof.* First, we fix  $\epsilon > 0$ . Thm. C.1 tells us that there is an  $N$  such that for all  $n > N$

$$\gamma - \epsilon < L_n^{1/n} < \gamma + \epsilon.$$

Combining the above equation with (C2), it is clear that if we can show that for all  $m > M$  we have

$$rm^\eta < (\gamma - \epsilon)^n < L_n < (\gamma + \epsilon)^n < sm, \quad (\text{C3})$$

then we have proven our theorem.

We note that  $n > N$  must be satisfied if we choose  $m$  large enough so that if the first inequality is satisfied, i.e.

$$n > \frac{\log r + \eta \log m}{\log(\gamma - \epsilon)} > N.$$

Additionally, when

$$\log s - \log r + (1 - \eta) \log m > 2 \log \gamma > 2 \log(\gamma - \epsilon),$$

there are at least two points  $n$  that satisfy

$$\log r + \eta \log m < n \log(\gamma - \epsilon) < \log s + \log m.$$

In particular, there is a point before the midpoint of the bound such that

$$\begin{aligned} \log r + \eta \log m &< n \log(\gamma - \epsilon) < \\ &\frac{1}{2}(\log r + \log s) + \frac{1 + \eta}{2} \log m. \end{aligned}$$

Using the above expression, we additionally know that

$$\begin{aligned} n \log(\gamma + \epsilon) &< \\ &\frac{\log(\gamma + \epsilon)}{\log(\gamma - \epsilon)} \left[ \frac{1}{2}(\log r + \log s) + \frac{1 + \eta}{2} \log m \right]. \end{aligned}$$

If we choose  $\epsilon$  small enough so that

$$\frac{\log(\gamma + \epsilon)}{\log(\gamma - \epsilon)} = \frac{2C}{1 + \eta},$$

where  $C < 1$ , then the above inequality reduces to

$$n \log(\gamma + \epsilon) < \frac{C}{1 + \eta}(\log r + \log s) + C \log m.$$

From here, it is clear that there exists an  $M$  so that  $m > M$  implies that

$$n \log(\gamma + \epsilon) < \log s + \log m,$$

giving our upper bound.  $\square$

## 2. The Polynomial

To prove Thm. III.2, we create a sequence of filters that annihilate a certain number of the frequencies  $\lambda_n$ . If we let  $0 < \alpha < 1/4$  and  $n \geq 1$ , we do this by defining the reference polynomial for frequency  $\omega$  and island period  $p$  as

$$\begin{aligned} q_{\alpha, n}(z) &= \prod_{j=1}^{\lfloor \alpha p L_n \rfloor} \frac{(z - \lambda_j)(z - \lambda_{-j})}{(1 - \lambda_j)(1 - \lambda_{-j})} \\ &\quad \times \prod_{j=\lfloor \alpha p L_n \rfloor + 1}^{\lfloor p L_n / 2 \rfloor} \frac{(z - \mu_j)(z - \mu_{-j})}{(1 - \mu_j)(1 - \mu_{-j})} \quad (\text{C4}) \end{aligned}$$

where for  $0 \leq \ell < p$  and  $j \geq 0$

$$\mu_{jp+\ell} = e^{2\pi i(jN_n/L_n + \ell)/p}.$$

A fraction  $2\alpha$  of the roots of  $q_{\alpha, n}$  exactly match  $\lambda_j$ , and the rest of the roots  $\mu_j$  lie on roots of unity. The polynomial has been chosen so that it is comparable to the polynomial with roots at roots of unity, which is well understood as a cardinal function of the discrete Fourier transform.

In particular, we have:

**Lemma C.3.** *Let  $0 < \omega \leq 1$  have a sequence of continued fractions approximants with  $\{N_n/L_n\}$ . For any  $n \geq 1$  and  $\alpha < 1/4$ , the reference polynomial obeys the bound*

$$|q_{\alpha,n}(z)| \leq C_\alpha$$

for  $|z| \leq 1$  and  $C_\alpha > 1$ , where  $C_\alpha$  is independent of  $\omega$  and  $n$ .

*Proof.* The proof consists of comparing the polynomial above to the polynomial with equispaced nodes

$$Q_{\alpha,n}(z) = \prod_{j=1}^{\lfloor L_n/2 \rfloor} \frac{(z - \mu_j)(z - \mu_{-j})}{(1 - \mu_j)(1 - \mu_{-j})}.$$

We begin by assuming  $L_n$  is odd. In this case, the polynomial coefficients of  $Q_{\alpha,n}$  can be thought of as a column of the inverse discrete Fourier transform. This is because we have  $Q(1) = 1$  and  $Q(\mu_j) = 0$  for  $1 \leq |j| \leq \lfloor L_n/2 \rfloor$ . So,  $Q_{\alpha,n}$  is simply the interpolation of a unit vector that is nonzero at 1, i.e.  $\mathbf{d} = F^{-1}\mathbf{e}_0$  where  $\mathbf{d}$  are the polynomial coefficients,  $F_{mn} = \mu_{n-M_n}^m$  is the orthogonal discrete Fourier transform mapping from Fourier coefficients to values at equispaced nodes, and  $M_n = \lfloor pL_n/2 \rfloor$ .

The problem of how Fourier interpolation changes when the nodes are perturbed is addressed in the paper by Yu and Townsend<sup>31</sup>. The authors prove an analogue to the Kadec-1/4 theorem<sup>32</sup> from sampling theory, which states that as long as the locations of Fourier nodes does not move a length farther than  $\alpha/4$  from the equispaced grid, the non-uniform Fourier transform matrix  $\tilde{F}$  does not differ too much in norm from the equispaced Fourier transform, i.e.

$$\|\tilde{F}\| \leq (1 + \phi_\alpha) \|F\|_2,$$

where  $\|F\| = L_n$ ,  $\phi_\alpha = 1 - \cos \pi\alpha + \sin \pi\alpha$ , and

$$\tilde{F}_{mn} = \begin{cases} \lambda_{n-M_n}, & |n - M_n| \leq \lfloor \alpha p L_n \rfloor, \\ \mu_{n-M_n}^m, & \lfloor \alpha p L_n \rfloor < |n - M_n| \leq M_n. \end{cases}$$

Additionally, via Weyl's inequality we have

$$\|\tilde{F}^{-1}\|_2 \leq \frac{1}{(1 - \phi_\alpha) \|F\|_2}.$$

We can apply these identities to our problem. Using the approximation bound (C1), we find that  $|\lambda_j - \mu_j| < \alpha\pi/L_n$  for  $|j| < \lfloor \alpha L_n \rfloor$ , obeying the hypothesis for the 1/4 theorem. Hence, the coefficients of  $q_{\alpha,n}$

$$\|\mathbf{c}_{\alpha,n}\|_2 = \|\tilde{F}^{-1}\mathbf{e}_0\|_2 \leq \frac{1}{(1 - \phi_\alpha)\sqrt{L_n}}.$$

The polynomial evaluation bound directly follows from the above inequality.  $\square$

### 3. Convergence of the Least-Squares Residual

With Lemma C.3 and Corollary C.2, we have enough information to prove Theorem III.2:

*Proof of Theorem III.2.* Choose  $r = s = 1/p$  and  $\eta < 1$ . Then, from Corollary C.2 we know there is an  $L$  such that  $2K + 1 > L$  implies there is an  $n$  such that the continued fraction representation of  $\omega$  obeys

$$(2K + 1)^\eta < pL_n < 2K + 1.$$

Therefore, the filter  $\mathbf{c}_{\alpha,n}$  with the coefficients of  $q_{\alpha,n}$  has a length less than the maximum allowed for the linear system. To make the filter compatible with the full system, we pad  $\mathbf{c}_{\alpha,n}$  symmetrically with zeros on the top and bottom. We note that  $\mathbf{c}_{\alpha,n}$  obeys the palindromic and correct-mean constraints by construction.

Now, consider the application of  $\mathbf{c}_{\alpha,n}$  on  $U$ . Using the Fourier representation of  $\mathbf{g}$ , we have

$$U_{jk} = \mathbf{g}((k + j)\omega) = \sum_{m \in \mathbb{Z}} \mathbf{g}_m \lambda_m^{k+j}.$$

Applying the filter, we find

$$\begin{aligned} (U\mathbf{c}_{\alpha,n})_j &= \sum_{m \in \mathbb{Z}} \mathbf{g}_m \lambda_m^j q_{\alpha,n}(\lambda_m), \\ &= \sum_{|m| > \lfloor \alpha p L_n \rfloor} \mathbf{g}_m \lambda_m^j q_{\alpha,n}(\lambda_m). \end{aligned}$$

Using Lemma C.3, we have

$$\begin{aligned} |(U\mathbf{c}_{\alpha,n})_j| &\leq C_\alpha \sum_{|m| > \lfloor \alpha p L_n \rfloor} \|\mathbf{g}_m\|, \\ &\leq C_\alpha C_M (\alpha p L_n)^{-M+1}, \\ &\leq C_\alpha C_M \alpha^{-M+1} (2K + 1)^{\eta(-M+1)}. \end{aligned}$$

where we use the fact that  $\mathbf{g} \in C^M$  for the second inequality for some  $C_M > 0$ . Substituting into the full RRE quadratic form, we have

$$\begin{aligned} (U\mathbf{c}_{\alpha,n})^T W_T (U\mathbf{c}_{\alpha,n}) &= \sum_{t=0}^T w_{t,T} (U\mathbf{c}_{\alpha,n})_j^2, \\ &\leq C(2K + 1)^{2\eta(-M+1)}. \end{aligned}$$

In the case that  $\mathbf{h}$  is real-analytic, there exists a constant  $\rho$  such that

$$\|\mathbf{g}_m\| < C_\rho \rho^{-m}.$$

This further implies that for some  $\tilde{C} > 0$

$$\begin{aligned} |(U\mathbf{c}_{\alpha,n})_j| &< \tilde{C} \rho^{-\alpha p L_n}, \\ &< \tilde{C} \rho^{-\alpha(2K+1)^\eta} \end{aligned}$$

After substitution into the quadratic form, we find our result.  $\square$

- <sup>1</sup>E. Kolemen, N. J. Kasdin, and P. Gurfil, “Multiple Poincaré sections method for finding the quasiperiodic orbits of the restricted three body problem,” *Celestial Mechanics and Dynamical Astronomy* **112**, 47–74 (2012).
- <sup>2</sup>E. J. Paul, S. R. Hudson, and P. Helander, “Heat conduction in an irregular magnetic field. Part 2. Heat transport as a measure of the effective non-integrable volume,” *Journal of Plasma Physics* **88**, 905880107 (2022), publisher: Cambridge University Press.
- <sup>3</sup>R. d. I. LLave, *A tutorial on KAM theory* (American Mathematical Society; Oxford University Press, Providence, R.I., Oxford, 2004).
- <sup>4</sup>A. Haro, M. Canadell, J.-L. Figueras, A. Luque, and J. M. Mondelo, *The Parameterization Method for Invariant Manifolds: From Rigorous Results to Effective Computations*, Applied Mathematical Sciences, Vol. 195 (Springer International Publishing, Cham, 2016).
- <sup>5</sup>J.-L. Figueras, A. Haro, and A. Luque, “Rigorous Computer-Assisted Application of KAM Theory: A Modern Approach,” *Foundations of Computational Mathematics* **17**, 1123–1193 (2017).
- <sup>6</sup>R. L. Dewar and J. D. Meiss, “Flux-minimizing curves for reversible area-preserving maps,” *Physica D: Nonlinear Phenomena* **57**, 476–506 (1992).
- <sup>7</sup>R. L. Dewar, S. R. Hudson, and P. F. Price, “Almost invariant manifolds for divergence-free fields,” *Physics Letters A* **194**, 49–56 (1994).
- <sup>8</sup>J. Guckenheimer and P. Holmes, *Nonlinear oscillations, dynamical systems, and bifurcations of vector fields*, corr. 5th print ed., Applied mathematical sciences No. v. 42 (Springer, New York, 1997).
- <sup>9</sup>S. Das and J. A. Yorke, “Super convergence of ergodic averages for quasiperiodic orbits,” *Nonlinearity* **31**, 491–501 (2018).
- <sup>10</sup>E. Sander and J. Meiss, “Birkhoff averages and rotational invariant circles for area-preserving maps,” *Physica D: Nonlinear Phenomena* **411**, 132569 (2020).
- <sup>11</sup>E. Sander and J. D. Meiss, “Rotation Vectors for Torus Maps by the Weighted Birkhoff Average,” (2023), arXiv:2310.11600 [nlin].
- <sup>12</sup>A. Luque and J. Villanueva, “Quasi-Periodic Frequency Analysis Using Averaging-Extrapolation Methods,” *SIAM Journal on Applied Dynamical Systems* **13**, 1–46 (2014).
- <sup>13</sup>J. Villanueva, “A new averaging-extrapolation method for quasiperiodic frequency refinement,” *Physica D: Nonlinear Phenomena* **438**, 133344 (2022).
- <sup>14</sup>D. Blessing and J. D. M. James, “Weighted Birkhoff Averages and the Parameterization Method,” (2023), arXiv:2306.16597 [math].
- <sup>15</sup>J. Laskar, “Introduction to Frequency Map Analysis,” in *Hamiltonian Systems with Three or More Degrees of Freedom*, edited by C. Simó (Springer Netherlands, Dordrecht, 1999) pp. 134–150.
- <sup>16</sup>S. Das, C. B. Dock, Y. Saiki, M. Salgado-Flores, E. Sander, J. Wu, and J. A. Yorke, “Measuring quasiperiodicity,” *EPL (Europhysics Letters)* **114**, 40005 (2016).
- <sup>17</sup>S. Das, Y. Saiki, E. Sander, and J. A. Yorke, “Quantitative quasiperiodicity,” *Nonlinearity* **30**, 4111–4140 (2017).
- <sup>18</sup>S. Das, Y. Saiki, E. Sander, and J. A. Yorke, “Solving the Babylonian problem of quasiperiodic rotation rates,” *Discrete & Continuous Dynamical Systems - S* **12**, 2279–2305 (2019).
- <sup>19</sup>N. Baresi, Z. P. Olikara, and D. J. Scheeres, “Fully Numerical Methods for Continuing Families of Quasi-Periodic Invariant Tori in Astrodynamics,” *The Journal of the Astronautical Sciences* **65**, 157–182 (2018).
- <sup>20</sup>J. W. Burby, Q. Tang, and R. Maulik, “Fast neural Poincaré maps for toroidal magnetic fields,” *Plasma Physics and Controlled Fusion* **63**, 024001 (2021).
- <sup>21</sup>P. Jin, Z. Zhang, A. Zhu, Y. Tang, and G. E. Karniadakis, “SympNets: Intrinsic structure-preserving symplectic networks for identifying Hamiltonian systems,” *Neural Networks* **132**, 166–179 (2020).
- <sup>22</sup>K. Rath, C. G. Albert, B. Bisl, and U. von Toussaint, “Symplectic Gaussian process regression of maps in Hamiltonian systems,” *Chaos: An Interdisciplinary Journal of Nonlinear Science* **31**, 053121 (2021).
- <sup>23</sup>A. Sidi, *Vector Extrapolation Methods with Applications* (Society for Industrial and Applied Mathematics, Philadelphia, PA, 2017).
- <sup>24</sup>V. Mandelshtam, “FDM: the filter diagonalization method for data processing in NMR experiments,” *Progress in Nuclear Magnetic Resonance Spectroscopy* **38**, 159–196 (2001).
- <sup>25</sup>Y. Coudène, *Ergodic Theory and Dynamical Systems*, Universitext (Springer, London, 2016).
- <sup>26</sup>C. C. Paige and M. A. Saunders, “LSQR: An Algorithm for Sparse Linear Equations and Sparse Least Squares,” *ACM Transactions on Mathematical Software* **8**, 43–71 (1982).
- <sup>27</sup>I. J. Good, “THE COLLEAGUE MATRIX, A CHEBYSHEV ANALOGUE OF THE COMPANION MATRIX,” *The Quarterly Journal of Mathematics* **12**, 61–68 (1961).
- <sup>28</sup><https://github.com/maxeruth/SymplecticMapTools.jl>.
- <sup>29</sup>L. N. Trefethen, *Approximation theory and approximation practice*, Applied mathematics (Society for Industrial and Applied Mathematics, Philadelphia, 2013).
- <sup>30</sup>A. I. Khinchin and H. Eagle, *Continued fractions* (Dover Publications, Mineola, N.Y., 1997).
- <sup>31</sup>A. Yu and A. Townsend, “On the stability of unevenly spaced samples for interpolation and quadrature,” *BIT Numerical Mathematics* **63**, 23 (2023).
- <sup>32</sup>M. I. Kadets, “The exact value of the paley-wiener constant,” in *Dokl. Akad. Nauk SSSR*, Vol. 155 (1964) pp. 1253–1254.