

## A. Mathematical Proof

### A.1. Proof of Thm. 5

**Lemma 3.** Consider  $\dot{x} = f(x)$  and  $\dot{y} = f(y)$ , where  $f$  is  $L$ -Lipschitz continuous. Then

$$\|x(t) - y(t)\| \leq \exp(Lt)\|x(0) - y(0)\|. \quad (10)$$

*Proof.* Since

$$\begin{aligned} x(t) &= x(0) + \int_0^t f(x(\tau))d\tau \\ y(t) &= y(0) + \int_0^t f(y(\tau))d\tau, \end{aligned}$$

triangular inequality and Lipschitz continuity give

$$\begin{aligned} \|x(t) - y(t)\| &\leq \|x(0) - y(0)\| + \int_0^t \|f(x(\tau)) - f(y(\tau))\|d\tau \\ &\leq \|x(0) - y(0)\| + \int_0^t L\|x(\tau) - y(\tau)\|d\tau. \end{aligned}$$

Gronwall inequality thus gives Eq. (10).  $\square$

**Lemma 4.** Consider  $\dot{x} = f(x)$ , with  $x(0) = x_0$  and  $L$ -Lipschitz continuous  $f$ . Then

$$\|x(h) - x_0\| \leq \frac{\exp(Lh) - 1}{L} \|f(x_0)\|. \quad (11)$$

(Note  $\frac{\exp(Lh) - 1}{L} = \mathcal{O}(h)$ .)

*Proof.* Note

$$\begin{aligned} x(h) &= x_0 + \int_0^h f(x(\tau))d\tau \\ &= x_0 + \int_0^h f(x(\tau)) - f(x_0) + f(x_0)d\tau. \end{aligned}$$

Let  $D(t) := x(t) - x_0$ . Then triangular inequality and Lipschitz continuity of  $f$  give

$$D(h) \leq \int_0^h LD(\tau) + \|f(x_0)\|d\tau$$

Gronwall lemma thus yields

$$D(h) \leq \exp(Lh)D(0) + \frac{\exp(Lh) - 1}{L} \|f(x_0)\|.$$

Since  $D(0) = 0$ , Eq. (11) is proved.  $\square$

*Proof of Thm. 5.* Let  $E_n = \|x(nh) - x_n\|$  denote the prediction accuracy, and  $\phi_{x_0}^h$  be the  $h$ -time flow map of the latent dynamics, i.e.,  $\phi_{x_0}^h := x(h)$  where  $x(\cdot)$  satisfies  $\dot{x} = f(x)$  subject to  $x(0) = x_0$ . Then

$$x((n+1)h) - x_{n+1} = x((n+1)h) - \phi_{x_n}^h + \phi_{x_n}^h - x_{n+1},$$

and therefore

$$E_{n+1} \leq \|x((n+1)h) - \phi_{x_n}^h\| + \|\phi_{x_n}^h - x_{n+1}\|.$$

The first term is exactly  $\|\phi_{x(nh)}^h - \phi_{x_n}^h\|$ , and by Lemma. 3, it is bounded by

$$\|\phi_{x(nh)}^h - \phi_{x_n}^h\| \leq \exp(Lh)\|x(nh) - x_n\| = \exp(Lh)E_n.$$

For the second term, Taylor expansion gives

$$\phi_{x_n}^h = x_n + hf(x_n) + h^2/2f'(\phi_{x_n}^\xi)f(\phi_{x_n}^\xi)$$

for some  $\xi \in [0, h]$ , and therefore

$$\begin{aligned} \|\phi_{x_n}^h - x_{n+1}\| &= \|h(f(x_n) - \tilde{f}(x_n)) + h^2/2f'(\phi_{x_n}^\xi)f(\phi_{x_n}^\xi)\| \\ &\leq h\delta + h^2/2\|f'(\phi_{x_n}^\xi)\|\|f(\phi_{x_n}^\xi)\|. \end{aligned}$$

Note  $\|f'\| \leq L$  as  $f$  is  $C^1$  and  $L$ -Lipschitz. For the  $f(\phi_{x_n}^\xi)$  factor, note Lemma. 4 gives

$$\|\phi_{x_n}^\xi - x_n\| \leq \frac{\exp(L\xi) - 1}{L} \|f(x_n)\|,$$

and therefore

$$\begin{aligned} \|f(\phi_{x_n}^\xi)\| &= \|f(x_n) + f(\phi_{x_n}^\xi) - f(x_n)\| \\ &\leq \|f(x_n)\| + \|f(\phi_{x_n}^\xi) - f(x_n)\| \\ &\leq \|f(x_n)\| + L\|\phi_{x_n}^\xi - x_n\| \\ &\leq \exp(L\xi)\|f(x_n)\| \end{aligned}$$

Since  $0 \leq \xi \leq h$ ,  $\exp(L\xi)$  is bounded. Moreover,  $f(x_n)$  is bounded because  $f$  is Lipschitz and therefore continuous and  $x_n$  is assumed to be bounded. Therefore, there exists constant  $C$  such that

$$\|f'(\phi_{x_n}^\xi)\|\|f(\phi_{x_n}^\xi)\| \leq C$$

Summarizing both terms, we have

$$E_{n+1} \leq E_n \exp(Lh) + h\delta + Ch^2/2.$$

Mathematical induction thus gives

$$\begin{aligned} E_N &\leq E_0 \exp(Lh)^N \\ &\quad + \left( \exp(Lh)^{N-1} + \exp(Lh)^{N-2} + \cdots + 1 \right) (h\delta + Ch^2/2) \\ &= E_0 \exp(LT) + \frac{\exp(LT) - 1}{\exp(Lh) - 1} (h\delta + Ch^2/2) \\ &\leq E_0 \exp(LT) + \frac{\exp(LT) - 1}{Lh} (h\delta + Ch^2/2) \\ &= \frac{\exp(LT) - 1}{L} (\delta + Ch/2). \end{aligned}$$

$\square$

### A.2. Proof of Thm. 3

**Definition 5** (Diophantine condition). A frequency vector  $\omega = \{\omega_1, \omega_2, \dots, \omega_d\}$  satisfies Diophantine condition if and only there exists positive constants  $\gamma, \nu$  such that  $\omega$  satisfies  $(\gamma, \nu)$ -Diophantine condition.

**Definition 6**  $((\gamma, \nu)$ -Diophantine set). For a set  $\Omega \subseteq \mathbb{R}^d$ , the corresponding  $(\gamma, \nu)$ -Diophantine set is defined as

$$\Omega^*(\gamma, \nu) \stackrel{\text{def}}{=} \left\{ \omega \in \Omega : \omega \text{ satisfy } (\gamma, \nu)\text{-Diophantine condition} \right\}.$$

**Definition 7** (Diophantine set). For a set  $\Omega \subseteq \mathbb{R}^d$ , the corresponding Diophantine set is defined as

$$\Omega^* \stackrel{\text{def}}{=} \bigcup_{\gamma > 0, \nu > 0} \Omega^*(\gamma, \nu).$$

**Theorem 6.** For any bounded domain  $\Omega \subseteq \mathbb{R}^d$ , there exists  $C > 0$ , such that the Lebesgue measure of the complementary of  $(\gamma, \nu)$ -Diophantine set with  $\nu \geq d$  is bounded from above,

$$\lambda(\Omega \setminus \Omega^*(\gamma, \nu)) \leq C \cdot \gamma.$$

*Proof.* See for instance (Hairer et al., 2006).  $\square$

**Theorem 7.** For any bounded domain  $\Omega \subseteq \mathbb{R}^d$ , Diophantine frequencies exist almost everywhere.

*Proof.* Since  $\lambda(\Omega \setminus \Omega^*) \leq \lambda(\Omega \setminus \Omega^*(\gamma, \nu))$ ,  $\forall \gamma > 0, \nu > 0$ , Thm. 6 gives,  $\forall \gamma > 0$ ,

$$\lambda(\Omega \setminus \Omega^*) \leq C \cdot \gamma,$$

meaning that Diophantine frequencies exist almost everywhere in  $\Omega$ .  $\square$

**Remark 7.** Even Diophantine frequencies exist almost everywhere in bounded domain  $\Omega \subseteq \mathbb{R}^d$ ,  $\Omega \setminus \Omega^*$  is still an open and dense set in  $\mathbb{R}^d$  (see for instance Hairer et al., 2006).

**Lemma 5** (Cauchy's inequality). Suppose that  $f$  is a holomorphic function on a closed ball  $\bar{\mathcal{B}}_r(\theta^*) \subset \mathbb{C}$  with  $r > 0$ . If  $|f(\theta)| \leq M$  for all  $\theta$  on the boundary of  $\mathcal{B}_r(\theta^*)$ , then for all  $n \geq 0$ ,

$$\left| f^{(n)}(\theta^*) \right| \leq \frac{n!M}{r^n}.$$

*Proof.* See for instance (Stein & Shakarchi, 2010).  $\square$

**Definition 8** (average over angles). Assume  $F(\theta)$  is periodic in each argument, i.e.,  $F: \mathbb{T}^d \rightarrow \mathbb{R}$ , then the (angle) average of  $F$  is defined as

$$\bar{F} = \frac{1}{(2\pi)^d} \int_{\mathbb{T}^d} F(\theta) d\theta. \quad (12)$$

**Definition 9** (complex extension of  $\mathbb{T}^d$ ). The complex extension of  $\mathbb{T}^d$  of width  $\rho$  is defined as

$$\mathcal{B}_\rho(\mathbb{T}^d) = \{ \theta \in \mathbb{T}^d + i\mathbb{R}^d; \|Im\theta\| < \rho \}. \quad (13)$$

**Definition 10.** For an analytic function  $f(\cdot) = [f_1(\cdot), \dots, f_d(\cdot)] \in \mathbb{C}^d$ , we define the following norm

$$\|f\|_{\infty, S} := \sum_{i=1}^d \sup_{x \in S} |f_i(x)|. \quad (14)$$

**Lemma 6.** Suppose  $\omega \in \mathbb{R}^d$  satisfies the  $(\gamma, \nu)$ -Diophantine condition and  $G(\theta) \in \mathbb{R}$  is a bounded and analytic function on  $\mathcal{B}_\rho(\mathbb{T}^d)$ . Then, with  $\bar{G}$  being the average of  $G(\theta)$ , the PDE

$$DF(\theta) \cdot \omega + G(\theta) = \bar{G} \quad (15)$$

has a unique real analytic solution  $F(\cdot)$  with  $\bar{F} = 0$ . Moreover, for every positive  $\delta < \min(\rho, 1)$ ,  $F$  is bounded on  $\mathcal{B}_{\rho-\delta}(\mathbb{T}^d)$  by

$$\begin{cases} \|F\|_{\infty, \mathcal{B}_{\rho-\delta}(\mathbb{T}^d)} \leq \kappa_0 \delta^{-\alpha+1} \|G\|_{\infty, \mathcal{B}_\rho(\mathbb{T}^d)}, \\ \|\partial_\theta F\|_{\infty, \mathcal{B}_{\rho-\delta}(\mathbb{T}^d)} \leq \kappa_1 \delta^{-\alpha} \|G\|_{\infty, \mathcal{B}_\rho(\mathbb{T}^d)}, \end{cases}$$

with  $\alpha = \nu + d + 1$  and  $\kappa_0 = \nu^{-1} 8^d 2^\nu \nu!$ ,  $\kappa_1 = \nu^{-1} 8^d 2^{\nu+1} (\nu + 1)!$ .

*Proof.* See for instance (Hairer et al., 2006).  $\square$

**Lemma 7.** Consider a nearly integrable system with generating function  $S(I_0, \varphi_1) = S_0(I_0) + \varepsilon S_1(I_0, \varphi_1)$ . Suppose  $S_0$  and  $S_1$  are analytic and bounded in a complex neighborhood of  $\mathcal{D}_1 \subseteq \mathbb{R}^d$  and  $\mathcal{D} = \mathcal{D}_1 \times \mathbb{T}^d$  respectively. Then, there exists a real analytic canonical transformation  $(J, \theta) \leftrightarrow (I, \varphi)$  generated by  $\mathcal{T}(J, \varphi) = J \cdot \varphi + \varepsilon \mathcal{T}_1(J, \varphi)$ , such that the generating function in  $J, \theta$  variables takes the form of

$$\tilde{S}(J_0, \theta_1) = \tilde{S}_0(J_0) + \varepsilon^2 \tilde{R}_2(J_0, \theta_1, \varepsilon). \quad (16)$$

with  $\varepsilon^2 \tilde{R}_2$  being a higher-order perturbation to a new integrable system  $\tilde{S}_0$ . Moreover, this result is constructive: the transformation  $J, \theta \leftrightarrow I, \varphi$  is given by  $\mathcal{T}$  through

$$\begin{cases} I_i = \partial_2 \mathcal{T}(J_i, \varphi_i), \\ \theta_i = \partial_1 \mathcal{T}(J_i, \varphi_i). \end{cases} \quad \forall i = 0, 1 \quad (17)$$

and  $\mathcal{T}_1$  is the solution to the PDE

$$\partial_2 \mathcal{T}_1(J, \varphi) \cdot \omega(J) + G_1(J, \varphi) = \bar{G}_1(J), \quad (18)$$

where  $\omega(\cdot) = \nabla S_0(\cdot)$ ,  $G_1(J, \varphi) = S_1(J, \varphi + h \nabla S_0(J))$ , and  $\bar{G}_1(J)$  is its angle average.

*Proof.* The generating function  $S$  in  $\mathbf{J}, \boldsymbol{\theta}$  variables ( $\tilde{S}$ ) can be converted in the following using Eq. (17),

$$\begin{aligned}
 \tilde{S}(\mathbf{J}_0, \boldsymbol{\theta}_1) &= S(\mathbf{I}_0, \boldsymbol{\varphi}_1) \\
 &= S(\mathbf{J}_0 + \varepsilon \partial_2 \mathcal{T}_1(\mathbf{J}_0, \boldsymbol{\varphi}_0), \boldsymbol{\varphi}_1) \\
 &= S_0(\mathbf{J}_0 + \varepsilon \partial_2 \mathcal{T}_1(\mathbf{J}_0, \boldsymbol{\varphi}_0)) \\
 &\quad + \varepsilon S_1(\mathbf{J}_0 + \varepsilon \partial_2 \mathcal{T}_1(\mathbf{J}_0, \boldsymbol{\varphi}_0), \boldsymbol{\varphi}_0 + h \partial_1 S(\mathbf{I}_0, \boldsymbol{\varphi}_1)) \\
 &\quad + \mathcal{O}(\varepsilon^2) \\
 &= S_0(\mathbf{J}_0) + \varepsilon \partial_2 \mathcal{T}_1(\mathbf{J}_0, \boldsymbol{\varphi}_0) \cdot \nabla S_0(\mathbf{J}_0) \\
 &\quad + \varepsilon S_1(\mathbf{J}_0, \boldsymbol{\varphi}_0 + h \nabla S(\mathbf{J}_0)) + \mathcal{O}(\varepsilon^2) \\
 &= S_0(\mathbf{J}_0) + \varepsilon \{ \underbrace{\partial_2 \mathcal{T}_1(\mathbf{J}_0, \boldsymbol{\varphi}_0) \cdot \nabla S_0(\mathbf{J}_0)}_{\text{angle independent}} \\
 &\quad + \underbrace{S_1(\mathbf{J}_0, \boldsymbol{\varphi}_0 + h \nabla S(\mathbf{J}_0))}_{\text{angle dependent}} \} + \mathcal{O}(\varepsilon^2)
 \end{aligned} \tag{19}$$

with  $h$  the constant in Eq. (4). Collecting all  $\mathcal{O}(\varepsilon^2)$  terms, denoting them by a remainder term  $\tilde{R}_2$ , and converting all angles in  $\tilde{R}_2$  to  $\boldsymbol{\theta}_1$ , we have

$$\begin{aligned}
 \tilde{S}(\mathbf{J}_0, \boldsymbol{\theta}_1) &= S_0(\mathbf{J}_0) + \varepsilon \{ \underbrace{\partial_2 \mathcal{T}_1(\mathbf{J}_0, \boldsymbol{\varphi}_0) \cdot \nabla S_0(\mathbf{J}_0)}_{\text{angle independent}} \\
 &\quad + \underbrace{S_1(\mathbf{J}_0, \boldsymbol{\varphi}_0 + h \nabla S(\mathbf{J}_0))}_{\text{angle dependent}} \} + \varepsilon^2 \tilde{R}(\mathbf{J}_0, \boldsymbol{\theta}_1, \varepsilon).
 \end{aligned} \tag{20}$$

As long as the terms underlined in Eq. (20) add up to a function of  $\mathbf{J}_0$  only,  $\tilde{S}(\mathbf{J}_0, \boldsymbol{\theta}_1)$  won't have angle dependence till the  $\mathcal{O}(\varepsilon^2)$  term. This leads to a solvability requirement. More precisely, let  $G_1(\mathbf{J}_0, \boldsymbol{\varphi}_0) = S_1(\mathbf{J}_0, \boldsymbol{\varphi}_0 + h \nabla S_0(\mathbf{J}_0))$  and  $\overline{G}_1(\mathbf{J}_0)$  be its angle average. Then the PDE

$$\partial_2 \mathcal{T}_1(\mathbf{J}_0, \boldsymbol{\varphi}_0) \cdot \nabla S_0(\mathbf{J}_0) + G_1(\mathbf{J}_0, \boldsymbol{\varphi}_0) = \overline{G}_1(\mathbf{J}_0) \tag{21}$$

has a solution  $\mathcal{T}_1$ , and it makes the underlined terms  $\overline{G}_1(\mathbf{J}_0)$ . Therefore,  $\mathcal{T}_1$  and hence the generating function  $\mathcal{T}$  can be solved for from Eq. (21). The generating function  $\tilde{S}$  in  $\mathbf{J}, \boldsymbol{\theta}$  variables takes the form

$$\tilde{S}(\mathbf{J}_0, \boldsymbol{\theta}_1) = \tilde{S}_0(\mathbf{J}_0) + \varepsilon^2 \tilde{R}_2(\mathbf{J}_0, \boldsymbol{\theta}_1, \varepsilon),$$

with  $\tilde{S}_0(\mathbf{J}_0) = S_0(\mathbf{J}_0) + \varepsilon \overline{G}_1(\mathbf{J}_0)$ .  $\square$

**Remark 8.** Note that boundedness of  $\tilde{R}_2$  requires some extra condition on  $\nabla S_0$  (being Diophantine at some point; see Lemma. 9 for details). With bounded  $\tilde{R}_2$ , the generating function  $S(\cdot, \cdot)$  in  $\mathbf{I}, \boldsymbol{\varphi}$  variables is near integrable of order  $\mathcal{O}(\varepsilon)$  while  $\tilde{S}(\cdot, \cdot)$  in  $\mathbf{J}, \boldsymbol{\theta}$  variables is near integrable of order  $\mathcal{O}(\varepsilon^2)$ . Therefore, under the transformation  $\mathcal{T}$ , we get a 'better' set of variables  $\mathbf{J}, \boldsymbol{\theta}$  instead of  $\mathbf{I}, \boldsymbol{\varphi}$ , as the  $\mathbf{J}, \boldsymbol{\theta}$  dynamics is closer to being integrable, hence the dynamics of  $\mathbf{J}, \boldsymbol{\theta}$  can be estimated for longer time.

**Remark 9.** As angles satisfy periodic boundary conditions,

$\mathcal{T}_1$  and  $S_1$  can be expanded in Fourier series

$$\begin{cases} \mathcal{T}_1(\mathbf{J}, \boldsymbol{\theta}) = \sum_{\mathbf{k} \in \mathbb{Z}^d} t_{\mathbf{k}}(\mathbf{J}) \cdot e^{i(\mathbf{k} \cdot \boldsymbol{\theta})}, \\ S_1(\mathbf{J}, \boldsymbol{\theta}) = \sum_{\mathbf{k} \in \mathbb{Z}^d} s_{\mathbf{k}}(\mathbf{J}) \cdot e^{i(\mathbf{k} \cdot \boldsymbol{\theta})}. \end{cases} \tag{22}$$

Plugging Eq. (22) into Eq. (18), we have

$$t_{\mathbf{k}} \cdot (\mathbf{k} \cdot \boldsymbol{\omega}(\mathbf{J})) + s_{\mathbf{k}} = 0.$$

Noting that if  $\boldsymbol{\omega}$  doesn't satisfy Diophantine condition,  $\mathbf{k} \cdot \boldsymbol{\omega}$  can be small and may even vanish for some  $\mathbf{k} \in \mathbb{Z}^d$ , meaning that under some circumstances, the transformation constructed by  $\mathcal{T}(\mathbf{J}, \boldsymbol{\varphi}) = \mathbf{J} \cdot \boldsymbol{\varphi} + \varepsilon \mathcal{T}_1(\mathbf{J}, \boldsymbol{\varphi})$  is no longer of near identity ( $\text{Id} + \mathcal{O}(\varepsilon)$ ) as  $\mathcal{T}_1$  is not of order  $\mathcal{O}(1)$  any more.

**Lemma 8.** Consider a nearly integrable system with generating function  $S(\mathbf{I}_0, \boldsymbol{\varphi}_1) = S_0(\mathbf{I}_0) + \varepsilon S_1(\mathbf{I}_0, \boldsymbol{\varphi}_1)$ . Suppose  $S_0$  and  $S_1$  are analytic and bounded in a complex neighborhood of  $\mathcal{D}_1 \subseteq \mathbb{R}^d$  and  $\mathcal{D} = \mathcal{D}_1 \times \mathbb{T}^d$  respectively. Then, there exists a real analytic canonical transformation  $(\mathbf{J}, \boldsymbol{\theta}) \leftrightarrow (\mathbf{I}, \boldsymbol{\varphi})$  generated by  $\mathcal{T}(\mathbf{J}, \boldsymbol{\varphi}) = \mathbf{J} \cdot \boldsymbol{\varphi} + \sum_{k=1}^{N-1} \varepsilon^k \cdot \mathcal{T}_k(\mathbf{J}, \boldsymbol{\varphi})$ , such that the dynamics produced by the original generating function  $S$  rewritten in  $\mathbf{J}, \boldsymbol{\theta}$  variables corresponds to a transformed generating function  $\tilde{S}$  given by

$$\tilde{S}(\mathbf{J}_0, \boldsymbol{\theta}_1) = \tilde{S}_0(\mathbf{J}_0) + \varepsilon^N \tilde{R}_N(\mathbf{J}_0, \boldsymbol{\theta}_1, \varepsilon), \tag{23}$$

where  $\varepsilon^N \tilde{R}_N$  is a high-order perturbation to a new integrable system  $\tilde{S}_0(\mathbf{J}_0)$ . Here,  $\mathbf{J}, \boldsymbol{\theta} \leftrightarrow \mathbf{I}, \boldsymbol{\varphi}$  is defined by  $\mathcal{T}$  through

$$\begin{cases} \mathbf{I}_i = \partial_2 \mathcal{T}(\mathbf{J}_i, \boldsymbol{\varphi}_i), \\ \boldsymbol{\theta}_i = \partial_1 \mathcal{T}(\mathbf{J}_i, \boldsymbol{\varphi}_i), \end{cases} \quad \forall i = 0, 1. \tag{24}$$

*Proof.* Apply  $\mathcal{T}(\mathbf{J}, \boldsymbol{\varphi}) = \mathbf{J} \cdot \boldsymbol{\varphi} + \varepsilon \mathcal{T}_1 + \varepsilon \mathcal{T}_2 + \dots + \varepsilon^{N-1} \mathcal{T}_{N-1}$  to  $\tilde{S}(\mathbf{J}_0, \boldsymbol{\theta}_1) = S(\mathbf{I}_0, \boldsymbol{\varphi}_1)$  like in the proof of Lemma. 7, Taylor expand, and put all  $\mathcal{O}(\varepsilon^N)$  terms into  $\tilde{R}_N$ . Then we have

$$\begin{aligned}
 \tilde{S}(\mathbf{J}_0, \boldsymbol{\theta}_1) &= S_0(\mathbf{J}_0) \\
 &\quad + \varepsilon (\partial_2 \mathcal{T}_1(\mathbf{J}_0, \boldsymbol{\varphi}_0) \cdot \boldsymbol{\omega}(\mathbf{J}_0) + G_1(\mathbf{J}_0, \boldsymbol{\varphi}_0)) \\
 &\quad + \varepsilon^2 (\partial_2 \mathcal{T}_2(\mathbf{J}_0, \boldsymbol{\varphi}_0) \cdot \boldsymbol{\omega}(\mathbf{J}_0) + G_2(\mathbf{J}_0, \boldsymbol{\varphi}_0)) \\
 &\quad + \dots \\
 &\quad + \varepsilon^N \tilde{R}_N(\mathbf{J}_0, \boldsymbol{\theta}_1, \varepsilon),
 \end{aligned} \tag{25}$$

for some functions  $G_1, \dots, G_{N-1}$  periodic in the angles, and  $\boldsymbol{\omega}(\cdot) := \nabla S(\cdot)$ . Similar to  $G_1$  in the proof of Lemma. 7,  $G_2, G_3, \dots$  can be explicitly computed from Taylor expansions, but our proof does not require their specific expressions. Lemma. 7 solved for  $\mathcal{T}_1$  (periodic in angles) by making the underlined expression independent of the angles.

Repeating a similar procedure at different orders of  $\varepsilon$ ,  $\mathcal{T}_i$  (periodic) can be obtained for  $i = 1, 2, \dots$ . Specifically,  $\mathcal{T}_i$  satisfies the PDE

$$\partial_2 \mathcal{T}_i(\mathbf{J}_0, \boldsymbol{\varphi}_0) \cdot \boldsymbol{\omega}(\mathbf{J}_0) + G_i(\mathbf{J}_0, \boldsymbol{\varphi}_0) = \overline{G}_i(\mathbf{J}_0) \quad (26)$$

(The existence of the solution to this is proved in Lemma. 6). In general,

$$\begin{aligned} \tilde{S}(\mathbf{J}_0, \boldsymbol{\theta}_1) &= S_0(\mathbf{J}_0) + \sum_{k=1}^{N-1} \varepsilon^k \overline{G}_k(\mathbf{J}_0) + \varepsilon^N \tilde{R}_N(\mathbf{J}_0, \boldsymbol{\theta}_1) \\ &= \tilde{S}_0(\mathbf{J}_0) + \varepsilon^N \tilde{R}_N(\mathbf{J}_0, \boldsymbol{\theta}_1), \end{aligned}$$

with  $\tilde{S}_0(\mathbf{J}_0) = S_0(\mathbf{J}_0) + \sum_{k=1}^{N-1} \varepsilon^k \overline{G}_k(\mathbf{J}_0)$ .  $\square$

Note that  $\tilde{R}_N$  is not necessarily uniformly bounded in the whole data domain in Lemmas. 7 and 8, but in most cases, the uniform boundedness can be established (Lemma. 9) and under that circumstance, we will be able to quantitatively estimate the  $\mathbf{J}, \boldsymbol{\theta}$  dynamics (Lemma. 10).

**Lemma 9.** *Consider a nearly integrable system with generating function  $S(\mathbf{I}_0, \boldsymbol{\varphi}_1) = S_0(\mathbf{I}_0) + \varepsilon S_1(\mathbf{I}_0, \boldsymbol{\varphi}_1)$ . Suppose  $S_0$  and  $S_1$  are analytic and bounded in a complex neighborhood of  $\mathcal{D}_1$  and  $\mathcal{D} = \mathcal{D}_1 \times \mathbb{T}^d \subseteq \mathbb{R}^d \times \mathbb{T}^d$  respectively. There exists a real analytic symplectic change of coordinates of order  $\mathcal{O}(\varepsilon)$ :  $(\mathbf{I}, \boldsymbol{\varphi}) \leftrightarrow (\mathbf{J}, \boldsymbol{\theta})$  and under this transformation, the generating function in  $\mathbf{J}, \boldsymbol{\theta}$  takes the form*

$$\tilde{S}(\mathbf{J}_0, \boldsymbol{\theta}_1) = \tilde{S}_0(\mathbf{J}_0) + \varepsilon^N \tilde{R}_N(\mathbf{J}_0, \boldsymbol{\theta}_1, \varepsilon).$$

Suppose that  $\omega(\mathbf{J}^*)$  satisfies the  $(\gamma, \nu)$ -Diophantine condition for some  $\mathbf{J}^* \in \mathcal{D}_1$ . Then, for any fixed  $N \geq 2$ , there exist positive constants  $\varepsilon_0, c, C, \rho$  such that if  $\varepsilon \leq \varepsilon_0$ , then

$$\left\| \tilde{R}_N(\cdot, \cdot) \right\|_{\infty, \overline{\mathcal{B}_{2\delta}(\mathbf{J}^*)} \times \overline{\mathcal{B}_\rho(\mathbb{T}^d)}} \leq C$$

with  $\delta = c(N^2 |\log \varepsilon|)^{-\nu-1}$ .

*Proof.* Applying the canonical transformation  $\mathcal{T}$  constructed in Lemma. 8,  $\exists \rho', C' > 0$  such that

$$\left\| \tilde{R}_N(\mathbf{J}^*, \cdot) \right\|_{\infty, \mathcal{B}_{\rho'}(\mathbb{T}^d)} \leq C'(N, d, \gamma, \nu) \quad (27)$$

Approximate  $S$  with respect to angle variables using Fourier series  $\hat{S}_m$  till term  $m \propto |\log \varepsilon|$  such that the error is of order  $\mathcal{O}(\varepsilon^N)$  in a complex neighborhood of the torus  $\{\mathbf{J} = \mathbf{J}^*, \boldsymbol{\varphi} \in \mathbb{T}^d\}$ . Since  $|\mathbf{k} \cdot \boldsymbol{\omega}(\mathbf{J}^*)| \geq \gamma \|\mathbf{k}\|_1^{-\nu}$ ,  $\forall \mathbf{k} \in \mathbb{Z}^d$ , then  $\exists$  sufficiently small  $c > 0$  such that

$$|\mathbf{k} \cdot \boldsymbol{\omega}(\mathbf{J})| \geq \frac{1}{2} \gamma \|\mathbf{k}\|_1^{-\nu}, \|\mathbf{k}\|_1 \leq Nm \quad (28)$$

for all  $\mathbf{J} \in \mathcal{B}_{2\delta}(\mathbf{J}^*)$  with  $\delta = c(N^2 |\log \varepsilon|)^{-\nu-1}$ . As the Fourier coefficients of  $\hat{S}_m$  vanishes for  $\|\mathbf{k}\|_1 > Nm$ , thus according to condition Eq. (28) and combining  $S = \hat{S}_m + \mathcal{O}(\varepsilon^N)$ ,  $\exists \rho'' > 0$  and  $C'' > 0$ , such that

$$\left\| \tilde{R}_N(\mathbf{J}, \cdot) \right\|_{\infty, \mathcal{B}_{\rho''}(\mathbb{T}^d)} \leq C''(N, d, \gamma, \nu) \quad (29)$$

for all  $\|\mathbf{J} - \mathbf{J}^*\| \leq 2\delta$ .

In general,  $\exists C$  and  $\varepsilon$  independent  $\rho$  such that

$$\left\| \tilde{R}_N(\cdot, \cdot) \right\|_{\infty, \overline{\mathcal{B}_{2\delta}(\mathbf{J}^*)} \times \overline{\mathcal{B}_\rho(\mathbb{T}^d)}} \leq C(N, d, \gamma, \nu).$$

(for the specific forms of  $C', C''$ , which are lengthy but obtainable using tools of Fourier series and Cauchy's inequality, see for instance (Hairer et al., 2006)).  $\square$

**Lemma 10.** *Consider a nearly integrable system with generating function  $S(\mathbf{I}_0, \boldsymbol{\varphi}_1) = S_0(\mathbf{I}_0) + \varepsilon S_1(\mathbf{I}_0, \boldsymbol{\varphi}_1)$ . Suppose  $S_0$  and  $S_1$  are analytic and bounded in a complex neighborhood of  $\mathcal{D}_1$  and  $\mathcal{D} = \mathcal{D}_1 \times \mathbb{T}^d \subseteq \mathbb{R}^d \times \mathbb{T}^d$  respectively. Then there exists a real analytic near identity symplectic change of coordinates  $(\mathbf{I}, \boldsymbol{\varphi}) \mapsto (\mathbf{J}, \boldsymbol{\theta})$  of order  $\mathcal{O}(\varepsilon)$  and under this transformation, the generating function  $\tilde{S}$  in  $\mathbf{J}, \boldsymbol{\theta}$  variables takes the form*

$$\tilde{S}(\mathbf{J}, \boldsymbol{\theta}) = \tilde{S}_0(\mathbf{J}) + \varepsilon^N \tilde{R}_N(\mathbf{J}, \boldsymbol{\theta}, \varepsilon).$$

where  $\tilde{S}_0$  only depends on actions. Suppose that  $\omega(\mathbf{J}^*)$  satisfies the  $(\gamma, \nu)$ -Diophantine condition for some  $\mathbf{J}^* \in \mathcal{D}_1$ . Then, for any fixed  $N \geq 2$ ,  $\exists$  positive constants  $\varepsilon_0, c, C, \rho$  such that if  $\varepsilon \leq \varepsilon_0$ , the dynamics of  $\mathbf{J}, \boldsymbol{\theta}$  (generated by  $\tilde{S}$ ) with  $\|\mathbf{J}_0 - \mathbf{J}^*\|_2 \leq c |\log \varepsilon|^{-\nu-1}$  satisfies

$$\begin{cases} \|\mathbf{J}_n - \mathbf{J}_0\|_2 \leq C n h \varepsilon^N, \\ \|\boldsymbol{\theta}_n - \tilde{\boldsymbol{\omega}}(\mathbf{J}_0) n h - \boldsymbol{\theta}_0\|_2 \\ \leq C \left( n^2 h^2 + n h |\log \varepsilon|^{\nu+1} \right) \varepsilon^N. \end{cases} \quad (30)$$

Here  $\omega(\cdot) = \nabla S_0(\cdot)$  and  $\tilde{\omega}(\cdot) = \nabla \tilde{S}_0(\cdot)$ .

*Proof.* According to Lemma. 9,  $\exists c > 0, \rho > 0, C' > 0$  such that for  $\delta = c(N^2 |\log \varepsilon|)^{-\nu-1}$ ,  $\mathbf{J} \in \overline{\mathcal{B}_\delta(\mathbf{J}^*)}$  and  $\boldsymbol{\theta} \in \mathcal{B}_\rho(\mathbb{T}^d)$ ,  $|\tilde{R}_N(\mathbf{J}, \boldsymbol{\theta})| \leq C'$ . As  $\forall \mathbf{J} \in \overline{\mathcal{B}_\delta(\mathbf{J}^*)}$ ,  $\overline{\mathcal{B}_\delta(\mathbf{J})} \subset \overline{\mathcal{B}_{2\delta}(\mathbf{J}^*)}$ ,  $|\tilde{R}_N(\tilde{\mathbf{J}}, \boldsymbol{\theta})| \leq C'$  for all  $\tilde{\mathbf{J}} \in \overline{\mathcal{B}_\delta(\mathbf{J})}$  and  $\boldsymbol{\theta} \in \mathcal{B}_\rho(\mathbb{T}^d)$ . Using Cauchy's inequality (Lemma. 5), we have

$$\left\| \partial_2 \tilde{R}_N \right\|_{\infty, \overline{\mathcal{B}_\delta(\mathbf{J}^*)} \times \overline{\mathcal{B}_\rho(\mathbb{T}^d)}} \leq C' \quad (31)$$

and

$$\left\| \partial_1 \tilde{R}_N \right\|_{\infty, \overline{\mathcal{B}_\delta(\mathbf{J}^*)} \times \overline{\mathcal{B}_\rho(\mathbb{T}^d)}} \leq \frac{C'}{\delta}. \quad (32)$$

Plug Eq. (31) in the dynamics of  $\mathbf{J}, \boldsymbol{\theta}$

$$\begin{cases} \mathbf{J}_i = \mathbf{J}_{i+1} + h\partial_2 \tilde{S}(\mathbf{J}_i, \boldsymbol{\theta}_{i+1}), \\ \boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i + h\partial_1 \tilde{S}(\mathbf{J}_i, \boldsymbol{\theta}_{i+1}), \end{cases} \quad (33)$$

we have

$$\begin{aligned} \|\mathbf{J}_{i+1} - \mathbf{J}_i\|_2 &\leq C' h \varepsilon^N, \quad \forall i \in \mathbb{N} \\ \implies \|\mathbf{J}_n - \mathbf{J}_0\|_2 &\leq C' n h \varepsilon^N. \end{aligned}$$

for the  $\mathbf{J}$  sequence. On the other hand, for  $\boldsymbol{\theta}$  sequence, plug Eq. (32) in Eq. (33), we have

$$\|\boldsymbol{\theta}_{i+1} - (\boldsymbol{\theta}_i + h\tilde{\omega}(\mathbf{J}_i))\|_2 \leq \frac{C'}{\delta} h \varepsilon^N. \quad (34)$$

Since  $\tilde{\omega}$  is analytic on a bounded domain,  $\tilde{\omega}$  is Lipschitz. Thus, changing  $\mathbf{J}_i$  in Eq. (34) to  $\mathbf{J}_0$ ,  $\exists C''$  such that

$$\|\boldsymbol{\theta}_{i+1} - (\boldsymbol{\theta}_i + h\tilde{\omega}(\mathbf{J}_0))\|_2 \leq C'' n h^2 \varepsilon^N + \frac{C'}{\delta} h \varepsilon^N.$$

Therefore, letting  $C = \max(C', C'')$ , we have

$$\begin{aligned} \|\boldsymbol{\theta}_n - (\boldsymbol{\theta}_0 + n h \tilde{\omega}(\mathbf{J}_0))\|_2 &\leq \sum_{i=0}^{n-1} \|\boldsymbol{\theta}_{i+1} - (\boldsymbol{\theta}_i + h\tilde{\omega}(\mathbf{J}_0))\|_2 \\ &\leq C n h \left( n h + \frac{1}{\delta} \right) \varepsilon^N \\ &\leq C n h \left( n h + |\log \varepsilon|^{\nu+1} \right) \varepsilon^N, \end{aligned}$$

and Eq. (30) is proved.  $\square$

*Proof of Thm. 3.* Since it is assumed that analytic  $S_h$  and  $S_h^\theta$  satisfy

$$\sum_{i=1,2} \|\partial_i S_h^\theta(\cdot, \cdot) - \partial_i S_h(\cdot, \cdot)\|_\infty \leq C_1 \varepsilon$$

on a bounded domain  $\mathcal{D}$ ,  $S_h^\theta$  is an  $\mathcal{O}(\varepsilon)$  perturbation of  $S_h$  (note they can also be different by an  $\mathcal{O}(1)$  constant, but adding a constant to a generating function does not change its induced dynamics, and we thus assume without loss of generality that there is no such constant difference). Therefore, as  $S_h$  is integrable,  $S_h^\theta$  can be written as  $S_h^\theta(\mathbf{I}_0, \boldsymbol{\varphi}_1) = S_h(\mathbf{I}_0) + \varepsilon S_1(\mathbf{I}_0, \boldsymbol{\varphi}_1)$  for some function  $S_1$  modeling the (normalized) perturbation, and is thus nearly integrable. The latent dynamics, i.e., the exact solution of the integrable  $S_h(\mathbf{I}_0)$  with initial condition  $\mathbf{I}_0, \boldsymbol{\varphi}_0$  is

$$\begin{cases} \mathbf{I}(t) = \mathbf{I}_0, \\ \boldsymbol{\varphi}(t) = (\boldsymbol{\varphi}_0 + \boldsymbol{\omega}(\mathbf{I}_0)t) \bmod 2\pi. \end{cases}$$

Applying Lemma. 10 with  $N \geq 3$  (so that the  $(nh)^2 \varepsilon^N$  term is of order  $\mathcal{O}(\varepsilon)$  when  $nh\varepsilon = \mathcal{O}(1)$ ), there exists a

near identity canonical transformation  $(\mathbf{I}, \boldsymbol{\varphi}) \mapsto (\mathbf{J}, \boldsymbol{\theta})$  of order  $\varepsilon$  such that the solution of  $S_h^\theta$  in  $\mathbf{J}, \boldsymbol{\theta}$  variable satisfies

$$\begin{cases} \|\mathbf{J}_n - \mathbf{J}_0\|_2 \leq C' \varepsilon, \\ \|\boldsymbol{\theta}_n - (\boldsymbol{\theta}_0 + \tilde{\omega}(\mathbf{J}_0) nh)\|_2 \leq C' \varepsilon n h, \end{cases} \quad (35)$$

for some constant  $C'$  ( $\varepsilon$  independent)  $\forall n \leq h^{-1} \varepsilon^{-1}$  (so that  $nh\varepsilon$  is of order  $\mathcal{O}(1)$ ) with  $\tilde{\omega}(\cdot)$  defined in Lemma. 10. Note that the canonical transformation holds for all  $(\mathbf{I}_i, \boldsymbol{\varphi}_i) \leftrightarrow (\mathbf{J}_i, \boldsymbol{\theta}_i)$ , we have  $\|\mathbf{I}_i - \mathbf{J}_i\|_2 \leq k\varepsilon$ ,  $\|\boldsymbol{\varphi}_i - \boldsymbol{\theta}_i\|_2 \leq k\varepsilon$ ,  $\forall i \in \mathbb{N}$  for some constant  $k > 0$  and  $\|\tilde{\omega}(\mathbf{J}_0) - \boldsymbol{\omega}(\mathbf{I}_0)\|_2 \leq k'\varepsilon$  for some positive constants  $k'$ . Applying triangular inequality,  $\exists C > 0$ , such that

$$\begin{cases} \|\mathbf{I}_n - \mathbf{I}_0\|_2 \leq C\varepsilon, \\ \|\boldsymbol{\varphi}_n - (\boldsymbol{\varphi}_0 + \boldsymbol{\omega}(\mathbf{I}_0) nh)\|_2 \leq C\varepsilon n h, \end{cases}$$

for  $n \leq h^{-1} \varepsilon^{-1}$ .  $\square$

**Remark 10.**  $h^{-1} \varepsilon^{-1}$  is actually a conservative bound for  $n$  and one can extend the bound to be  $h^{-1} \varepsilon^{(1-N)/2}$ ,  $\forall N \geq 3$ . Since  $N$  can be arbitrary, even if  $\varepsilon$  cannot be made infinitesimal, as long as it is below a threshold, the time of validity of the error bound in Thm. 3 can be extended to arbitrarily long.

## B. Experimental Details

Like most neural-network based algorithms for learning dynamics, the full potential of GFNN is achieved in the data rich regime. When preparing training data, we not only prepared in an unbiased way, but also emphasized on fair comparisons so that each of the existing methods is given the same or more training data.

More specifically, for each experiment, the training set contains a number of sequences starting with different initial conditions. When training for predicting Hamiltonian dynamics (continuous), each sequence in the training set is stroboscopically sampled from simulated ground truth, which is obtained using high-order numerical integrator with sufficiently small timestep  $\tau \ll h$ . For each experiment, the data set with sequences of length 2 will be denoted as  $\mathcal{D}_2$ , and the data set with sequences of length 5 will be denoted as  $\mathcal{D}_5$ . VFNN, HNN, SRNN (seq\_len=2), and GFNN are trained with the same data set  $\mathcal{D}_2$ , while SRNN (seq\_len=5) is trained with  $\mathcal{D}_5$ . Note the number of flow maps ( $\phi$ ) needed for each sequence in  $\mathcal{D}_5$  is 4, while the number of maps for each sequence of  $\mathcal{D}_2$  is 1. Therefore, for fairness, the number of sequences in  $\mathcal{D}_2$ ,  $n_{train}(\mathcal{D}_2)$ , is set to be four times  $n_{train}(\mathcal{D}_5)$  in most examples (exceptions will be explained).

All experiments are performed with PyTorch (CUDA) on a machine with GeForce RTX 3070 graphic card, AMD



Ryzen 7 3700X 8-Core Processor, 16 GB memory and the Linux distribution of openSUSE Leap 15.2.

Codes are provided.

### B.1. 2-Body Problem

The step size of each data sequence in  $\mathcal{D}_2$  and  $\mathcal{D}_5$  is  $h = 0.1$ . The ground truth trajectory is simulated using a 4th order symplectic integrator with step size  $10^{-3}$ . The initial condition of each data sequence is uniformly drawn from the orbits with semi-major axis  $a \in (0.8, 1.2)$ , eccentricity  $e \in (0, 0.05)$ . In terms of the number of samples,  $n_{train}(\mathcal{D}_2) = 100,000$ ,  $n_{train}(\mathcal{D}_5) = 100,000$ . Note SRNN (seq.len=5) is provided more training data  $n_{train}(\mathcal{D}_5)$  than described above, which would be 25,000 instead, because less training data didn't provide good performance. The time derivative data of the vector field based methods (VFNN, HNN) are generated using (1st-order) finite difference.

$S_h^\theta$  is represented using multilayer perceptron (MLP), with 5 layers and 200, 100, 50, 20 neurons in hidden layers. The Adam optimizer is utilized with batch size 200. The model is trained for more than 20 epochs with initial learning rate 0.01. HNN, SRNN, SympNets are trained by their provided codes. HNN, SRNN are trained under default training setups and SympNets is trained using LA-SympNets with 30 layers and 10 sublayers (deeper than their default setups for improved performance).

### B.2. Hénon-Heiles System

The step size of each data sequence in  $\mathcal{D}_2$  and  $\mathcal{D}_5$  is  $h = 0.5$ . The ground truth trajectory is simulated using a 4th order symplectic integrator with step size  $10^{-3}$ . The initial condition of each data sequence is drawn from a centered Gaussian perturbation of states along one orbit randomly with variance  $0.01^2$ . In terms of the number of samples,  $n_{train}(\mathcal{D}_2) = 100,000$ ,  $n_{train}(\mathcal{D}_5) = 25,000$ . The data sets for the regular motion experiment and for the chaotic dynamics experiment are generated separably around a trajectory with energy level  $\frac{1}{12}$  and  $\frac{1}{6}$  respectively. The time derivative data of the vector field based methods (VFNN, HNN) are generated using (1st-order) finite difference.

The MLP that represents  $S_h^\theta$  has 5 layers and 200, 100, 50, 20 neurons in hidden layers. The Adam optimizer is utilized with batch size 200. The model is trained for more than 20 epochs with initial learning rate 0.01. HNN, SRNN are trained by their provided codes under default training setups.

### B.3. PCR3BP

The step size of each data sequence in  $\mathcal{D}_2$  and  $\mathcal{D}_5$  is  $h = 0.1$ . The ground truth trajectory is simulated using RK4 with step size  $10^{-3}$ . The initial condition of each data sequence is

drawn from a centered Gaussian perturbation of states along one orbit randomly with variance  $0.05^2$ . In terms of the number of samples,  $n_{train}(\mathcal{D}_2) = 100,000$ ,  $n_{train}(\mathcal{D}_5) = 25,000$ . The time derivative data of the vector field based methods (VFNN, HNN) are generated using (1st-order) finite difference.

The MLP that represents  $S_h^\theta$  has 5 layers and 200, 100, 50, 20 neurons in hidden layers. The Adam optimizer is utilized with batch size 200. The model is trained for more than 20 epochs with initial learning rate 0.01. HNN, SRNN are trained by their provided codes under default training setups.

### B.4. Standard Map

The step size of each data sequence in  $\mathcal{D}_2$  and  $\mathcal{D}_5$  is  $h = 1$ . The ground truth map is directly evolved from the discrete-time evolution map Eq. (9). The initial condition of each data sequence is drawn from a Gaussian perturbation of states along one orbit randomly with variance  $0.5^2$ . In terms of the number of samples in training / testing data,  $n_{train}(\mathcal{D}_2) = 1,000,000$ ,  $n_{train}(\mathcal{D}_5) = 250,000$ . The data sets of the regular motion experiment and for the chaotic dynamics experiment are generated separably with  $K = 0.6$  and  $K = 1.2$  and correspondingly different initial conditions respectively. The time derivative data of the vector field based methods (VFNN, HNN) are generated using (1st-order) finite difference (with  $\Delta t = 1$ ).

The MLP that represents  $S_h^\theta$  has 5 layers and 500, 500, 200, 20 neurons in hidden layers. The Adam optimizer is utilized with batch size 1000. The model is trained for more than 20 epochs with initial learning rate 0.001. HNN, SRNN are trained by their provided codes under default training setups.