



HAL
open science

A formal series approach to averaging: exponentially small error estimates

Philippe Chartier, Ander Murua, Jesus Maria Sanz-Serna

► **To cite this version:**

Philippe Chartier, Ander Murua, Jesus Maria Sanz-Serna. A formal series approach to averaging: exponentially small error estimates. *Discrete and Continuous Dynamical Systems - Series A*, 2012, 32 (9), 10.3934/dcds.2012.32.3009 . hal-00777178

HAL Id: hal-00777178

<https://inria.hal.science/hal-00777178>

Submitted on 17 Jan 2013

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A formal series approach to averaging: exponentially small error estimates

Ph. Chartier,^{*} A. Murua,[†] J. M. Sanz-Serna[‡]

February 16, 2012

Abstract

The techniques, based on formal series and combinatorics, used nowadays to analyze numerical integrators may be applied to perform high-order averaging in oscillatory periodic or quasi-periodic dynamical systems. When this approach is employed, the averaged system may be written in terms of (i) scalar coefficients that are universal, i.e. independent of the system under consideration and (ii) basis functions that may be written in an explicit, systematic way in terms of the derivatives of the Fourier coefficients of the vector field being averaged. The coefficients may be recursively computed in a simple fashion. We show that this approach may be used to obtain exponentially small error estimates, as those first derived by Neishtadt. All the constants that feature in the estimates have a simple explicit expression.

1 Introduction

This paper continues the work in [7] and [8] on the relations between the method of averaging (see e.g. [13], [18], [1, Chapt. 4], [2, Chapt. 10]) and the formal series expansions that are nowadays used to analyze numerical integrators [12], [19], [14], [11]. We show here how the approach introduced in [7] and [8] may be readily applied to derive exponentially small error estimates similar to those first proved by Neishtadt [16].

Let us assume that the problem to be averaged has been rewritten [4], [5] to take the familiar format:

$$\frac{d}{dt}y = \epsilon f(y, t\omega), \quad (1)$$

$$y(0) = y_0 \in \mathbb{R}^D, \quad (2)$$

where ϵ is a small parameter, $f = f(y, \theta)$ is sufficiently smooth and 2π -periodic in each of the components θ^j , $j = 1, \dots, d$, of θ , i.e. $\theta \in \mathbb{T}^d$, and $\omega \in \mathbb{R}^d$ is a constant vector of angular

^{*}INRIA Rennes and ENS Cachan Bretagne, Campus Ker-Lann, av. Robert Schumann, 35170 Bruz, France. Email: Philippe.Chartier@inria.fr

[†]Konputazio Zientziak eta A. A. Saila, Informatika Fakultatea, UPV/EHU, E-20018 Donostia-San Sebastián, Spain. Email: Ander.Murua@ehu.es

[‡](Corresponding author) Departamento de Matemática Aplicada, Facultad de Ciencias, Universidad de Valladolid, Valladolid, Spain. Email: sanzsern@mac.uva.es

frequencies. We assume throughout that ω is *non-resonant*, i.e. that, for each multi-index $\mathbf{k} \in \mathbb{Z}^d$, $\mathbf{k} \cdot \omega \neq 0$ (resonant problems of this form may be recast as non-resonant by lowering the number d of frequencies). The problem (1)–(2) is to be studied in an interval $0 \leq t \leq L/\epsilon$ so that y undergoes variations of order $\mathcal{O}(1)$. When $d = 1$ the right-hand side of (1) is *periodic* in the variable t ; for $d > 1$ the time-dependence is *quasi-periodic*.

The method of averaging considers a change of variables

$$y = Y + \epsilon \check{U}(Y, t\omega, \epsilon)$$

with $\check{U}(Y, \theta, \epsilon)$ 2π -periodic in θ^j , $j = 1, \dots, d$, that transforms (1)–(2) into

$$\frac{d}{dt}Y = \epsilon (F(Y, \epsilon) + R(Y, t\omega, \epsilon)), \quad Y(0) = Y_0, \quad (3)$$

where Y_0 is defined implicitly by $y_0 = Y_0 + \epsilon \check{U}(Y_0, \mathbf{0}, \epsilon)$ and R is a small reminder, so that in the transformed system the explicit time-dependence of the right-hand side has been ‘almost’ eliminated.

In the simplest version, F is the average f_0 of $f(\cdot, \theta)$ over the torus \mathbb{T}^d and $R = \mathcal{O}(\epsilon)$; accordingly the truncated autonomous problem $(d/dt)Y = \epsilon F(Y)$, $Y(0) = Y_0$, describes, with an $\mathcal{O}(\epsilon)$ error, the $\mathcal{O}(1)$ changes in the solution on $0 \leq t \leq L/\epsilon$. In higher-order versions, $R = \mathcal{O}(\epsilon^N)$, $N = 2, 3, \dots$, so that, if $Y(t)$ is the solution of the truncated problem, then $Y(t) + \epsilon \check{U}(Y(t), t\omega, \epsilon)$ provides an approximation to the solution $y(t)$ of (1)–(2) with an error of size $\mathcal{O}(\epsilon^N)$.

Although there are many possible variants, it is often the case that the transformed system in (3) is of the form

$$\frac{d}{dt}Y = \epsilon (F_1(Y) + \dots + \epsilon^{N-1} F_N(Y) + R^{(N)}(Y, t\omega, \epsilon)), \quad R^{(N)} = \mathcal{O}(\epsilon^N), \quad (4)$$

where the F_n do not change with N . Typically, the F_n , $n = 1, 2, \dots$, are found recursively starting from $F_1(Y) = f_0(Y)$; once F_1, \dots, F_N have been found, one changes variables in (4) so as to reduce the explicit time-dependence in the right-hand side to size $\epsilon \mathcal{O}(\epsilon^{N+1})$ and in this way one obtains F_{N+1} and a new averaged system with a higher-order reminder $R^{(N+1)}$. By letting $N \uparrow \infty$ in such a procedure, one writes a series

$$\sum_{n=1}^{\infty} \epsilon^{n-1} F_n(Y) \quad (5)$$

that, if convergent with sum $F(Y, \epsilon)$, leads to an averaged system $(d/dt)Y = \epsilon F(Y, \epsilon)$ where the time-dependence of the right-hand side has been completely eliminated. However the series (4) typically diverges and the complete suppression of the time-dependence is impossible.

Neishtadt [16] proved that in the periodic ($d = 1$) case, if f depends analytically on y and continuously on θ , it is possible to obtain an averaged problem (3) where the reminder R is *exponentially small* with respect to ϵ . Such an averaged system is found by performing successively $\mathcal{O}(1/\epsilon)$ intermediate changes of variables; each of these changes halves the magnitude of the remainder. Neishtadt’s exponential bounds have a number of very important consequences. In particular they imply [16, Prop. 3] that, under suitable hypotheses, symplectic integrators preserve energy with a small error over periods of time that are exponentially

long in the step-length; this is one of the main results in symplectic integration [19], [12] and plays a crucial role in e.g. molecular dynamics simulations. The article [17] by Ramis and Schäfke derives exponential error bounds for the periodic case using a technique different from that in [16]. Simó [20] has extended Neishtadt analysis to the quasi-periodic ($d > 1$) scenario; in that case f is demanded to depend analytically on y and θ and furthermore the vector ω of angular frequencies has to satisfy a diophantine condition.

Described in [8] is a method that allows the construction of a series of the form (5) such that, for each $N = 1, 2, \dots$, the system (1) may be transformed into (4) by a suitable change of variables. When this technique is used, each F_n is expressed as a combination of two elements of two sorts:

1. Scalar constant *coefficients* that are *universal*, i.e. independent of the function f in (1), and can be computed once and for all by means of simple recursions.
2. *Basis functions* constructed in an explicit, systematic way in terms of the derivatives of the Fourier coefficients of f .

With this methodology the coefficients F_n may be found independently of the required change of variables $Y \mapsto y$. In the present paper we show how the techniques in [8] may be applied to derive Neishtadt's exponential bounds. Here we bound, via Cauchy estimates, the size of the F_n in (4) and then, for a given value of ϵ , we determine how to choose $N = N(\epsilon)$ to minimize the magnitude of $R^{(N(\epsilon))}$.¹ This procedure leads to an exponentially small error bound (see Theorem 3.4) where all required constants have a simple, explicit expression. In a forthcoming publication we shall extend the present work to cover the quasi-periodic case as in [20].

The paper contains four sections. Section 2 reviews the approach to averaging described in [8]. The presentation is different from that in [8] because here we focus on the representation of the F_n in terms of so-called *word-basis* functions, while [8] emphasizes a more general but less compact representation in terms of *elementary differentials*. Section 3 presents the exponentially small error bounds and Section 4 contains a number of auxiliary results.

2 Formal series expansions

2.1 The expansion of y

While, as outlined in the introduction, standard approaches to high-order averaging envisage successive changes of variables to weaken the time-dependence of the vector field f in (1), the starting point of the technique in [7], [8] is the analysis of the expansion in powers of ϵ of the solution of (1)–(2). In this subsection we derive that expansion with the help of Lie operators.

For the time being we assume that f in (1) is given by a Fourier expansion

$$f(y, \theta) = \sum_{\mathbf{k} \in \mathbb{Z}^d} e^{i\mathbf{k} \cdot \theta} f_{\mathbf{k}}(y), \quad (6)$$

where the coefficients $f_{\mathbf{k}}$ are indefinitely continuously differentiable and, except for a finite number of values of $\mathbf{k} \in \mathbb{Z}^d$, vanish identically. We initially place stringent hypotheses on f

¹It is also possible to work as in [17] in order to avoid truncation of the formal series at the smallest term.

so as not to clutter the presentation with unwelcome details; as we show later (end of Section 2.4) it is also possible to work under much weaker alternative assumptions.

With each coefficient $f_{\mathbf{k}}$, $\mathbf{k} \in \mathbb{Z}^d$, we associate a linear differential operator

$$E_{\mathbf{k}} := \sum_{i=1}^D f_{\mathbf{k}}^i \frac{\partial}{\partial y^i}$$

(superscripts correspond to components). The operator $E_{\mathbf{k}}$ acting on a smooth function $g : \mathbb{R}^D \rightarrow \mathbb{R}$ yields a new smooth function $E_{\mathbf{k}}[g]$ defined by

$$\forall y \in \mathbb{R}^D, \quad E_{\mathbf{k}}[g](y) = \sum_{j=1}^D f_{\mathbf{k}}^j(y) \frac{\partial}{\partial y^j} g(y),$$

or in a more compact form, $E_{\mathbf{k}}[g](y) = \partial_y g(y) f_{\mathbf{k}}(y)$. Let us also consider the one-parameter family of linear operators Φ_t acting on smooth functions g such that $\Phi_t[g](y(0)) = g(y(t))$ for each solution $y(t)$ of (1) and each smooth function $g : \mathbb{R}^D \rightarrow \mathbb{R}$. Obviously, Φ_0 is the identity operator that we denote by I . We may write

$$\frac{d}{dt} g(y(t)) = \epsilon \partial_y g(y(t)) f(y(t), t\omega) = \epsilon \sum_{\mathbf{k} \in \mathbb{Z}^d} e^{i\mathbf{k} \cdot \omega t} E_{\mathbf{k}}[g](y(t)),$$

or equivalently

$$\frac{d}{dt} \Phi_t[g](y(0)) = \epsilon \sum_{\mathbf{k} \in \mathbb{Z}^d} e^{i\mathbf{k} \cdot \omega t} E_{\mathbf{k}}[g](y(0)).$$

This shows that Φ_t can be seen as the solution of the initial value problem

$$\frac{d}{dt} \Phi_t = \epsilon \sum_{\mathbf{k} \in \mathbb{Z}^d} e^{i\mathbf{k} \cdot \omega t} \Phi_t E_{\mathbf{k}}, \quad \Phi_0 = I,$$

that when solved by Picard iteration leads to the expansion

$$\Phi_t = I + \sum_{n=1}^{\infty} \epsilon^n \sum_{\mathbf{k}_1, \dots, \mathbf{k}_n \in \mathbb{Z}^d} \alpha_{\mathbf{k}_1 \dots \mathbf{k}_n}(t) E_{\mathbf{k}_1} \cdots E_{\mathbf{k}_n},$$

where

$$\alpha_{\mathbf{k}_1 \dots \mathbf{k}_n}(t) := \int_0^t e^{i\mathbf{k}_n \cdot \omega t_n} dt_n \int_0^{t_n} e^{i\mathbf{k}_{n-1} \cdot \omega t_{n-1}} dt_{n-1} \cdots \int_0^{t_2} e^{i\mathbf{k}_1 \cdot \omega t_1} dt_1. \quad (7)$$

By definition of Φ_t , for each smooth g and each solution $y(t)$ of (1)

$$g(y(t)) = g(y(0)) + \sum_{n=1}^{\infty} \epsilon^n \sum_{\mathbf{k}_1, \dots, \mathbf{k}_n \in \mathbb{Z}^d} \alpha_{\mathbf{k}_1 \dots \mathbf{k}_n}(t) E_{\mathbf{k}_1} \cdots E_{\mathbf{k}_n}[g](y(0)),$$

an expansion that, with the complex exponentials replaced by arbitrary scalar functions $u_k(t)$ (the controls), is known in non-linear control theory as the Chen-Fliess series [9], [10], [21].

The values $y(t)$ can be recovered by considering this series for each of the coordinate functions $g_i(y) = y^i$, $i = 1, \dots, D$. In this way we conclude that the solution $y(t)$ of (1)–(2) can be formally expanded as

$$y(t) = y_0 + \sum_{n=1}^{\infty} \epsilon^n \sum_{\mathbf{k}_1, \dots, \mathbf{k}_n \in \mathbb{Z}^d} \alpha_{\mathbf{k}_1 \dots \mathbf{k}_n}(t) f_{\mathbf{k}_1 \dots \mathbf{k}_n}(y_0), \quad (8)$$

with

$$f_{\mathbf{k}_1 \dots \mathbf{k}_n}(y) := \partial_y f_{\mathbf{k}_2 \dots \mathbf{k}_n}(y) f_{\mathbf{k}_1}(y). \quad (9)$$

The coefficient

$$\sum_{\mathbf{k}_1, \dots, \mathbf{k}_n \in \mathbb{Z}^d} \alpha_{\mathbf{k}_1 \dots \mathbf{k}_n}(t) f_{\mathbf{k}_1 \dots \mathbf{k}_n}(y_0) \quad (10)$$

of ϵ^n in (8) is made up of two ingredients: the *word-basis functions*² $f_{\mathbf{k}_1 \dots \mathbf{k}_n}(y)$ defined in (9) and the scalar *coefficients* $\alpha_{\mathbf{k}_1 \dots \mathbf{k}_n}$ in (7). Of importance is the fact that the coefficients depend on the frequencies ω but are completely independent of the function $f(\cdot, \cdot)$ in (6). On the other hand, the word-basis functions only depend on the Fourier coefficients $f_{\mathbf{k}}$; more precisely they are combinations of partial derivatives of these coefficients, e.g.

$$f_{\mathbf{1}\mathbf{k}}(y) = \partial_y f_{\mathbf{k}}(y) f_{\mathbf{1}}(y), \quad (11)$$

$$\begin{aligned} f_{\mathbf{m}\mathbf{l}\mathbf{k}}(y) &= \partial_y f_{\mathbf{l}\mathbf{k}}(y) f_{\mathbf{m}}(y) \\ &= \partial_{yy} f_{\mathbf{k}}(y) [f_{\mathbf{l}}(y), f_{\mathbf{m}}(y)] + \partial_y f_{\mathbf{k}}(y) \partial_y f_{\mathbf{l}}(y) f_{\mathbf{m}}(y). \end{aligned} \quad (12)$$

Here $\partial_y f_{\mathbf{k}}(y)$ is the first-order Fréchet derivative (Jacobian matrix) of $f_{\mathbf{k}}$ evaluated at y and $\partial_y f_{\mathbf{k}}(y) f_{\mathbf{l}}(y)$ is a matrix/vector product. Similarly $\partial_{yy} f_{\mathbf{k}}(y)$ is the second-order Fréchet derivative of $f_{\mathbf{k}}$ at y and $\partial_{yy} f_{\mathbf{k}}(y) [f_{\mathbf{l}}(y), f_{\mathbf{m}}(y)]$ denotes its action on the D -vectors $f_{\mathbf{l}}(y)$, $f_{\mathbf{m}}(y)$; thus the j -th component of $\partial_{yy} f_{\mathbf{k}}(y) [f_{\mathbf{m}}(y), f_{\mathbf{l}}(y)]$ is

$$\sum_{k, \ell} \left(\frac{\partial^2}{\partial y^k \partial y^\ell} f_{\mathbf{k}}^j(y) \right) f_{\mathbf{l}}^k(y) f_{\mathbf{m}}^\ell(y).$$

The functions $\partial_y f_{\mathbf{k}}(y) f_{\mathbf{l}}(y)$ in (11) and $\partial_{yy} f_{\mathbf{k}}(y) [f_{\mathbf{l}}(y), f_{\mathbf{m}}(y)]$, $\partial_y f_{\mathbf{k}}(y) \partial_y f_{\mathbf{l}}(y) f_{\mathbf{m}}(y)$ in (12) are instances of so-called *elementary differentials* (relative to the function f in (1)), see [8, Sec. 2.2]. The elementary differential $\partial_y f_{\mathbf{k}} f_{\mathbf{l}}$ is of order 2 (it contains two factors) and $\partial_{yy} f_{\mathbf{k}} [f_{\mathbf{l}}, f_{\mathbf{m}}]$, $\partial_y f_{\mathbf{k}} \partial_y f_{\mathbf{l}} f_{\mathbf{m}}$ have both order 3. In general, as illustrated by (11)–(12), each word-basis function $f_{\mathbf{k}_1 \dots \mathbf{k}_n}$ is a linear combination of elementary differentials of order n . In this way elementary differentials may be conceived of as building blocks to construct word-basis functions. A given elementary differential may enter in different word-basis function: thus $\partial_{yy} f_{\mathbf{k}} [f_{\mathbf{l}}, f_{\mathbf{m}}]$, $\mathbf{l} \neq \mathbf{m}$ enters in both $f_{\mathbf{m}\mathbf{l}\mathbf{k}}$ and $f_{\mathbf{l}\mathbf{m}\mathbf{k}}$. The elementary differentials of order n may be indexed by graphs called (mode-colored rooted) trees with n vertices (see [8, Table 1],). Since the work of J. Butcher in the 1960's [11], [19], [11], indexing the terms of a formal expansion by means of trees and other graphs has been standard in the analysis of numerical integrators.

²The terminology ‘word-basis’ will be motivated below.

2.2 Word-series

The notation in (10) and in expressions that appear later may be simplified by considering words $\mathbf{k}_1 \cdots \mathbf{k}_n$, made of letters \mathbf{k}_r , $r = 1, \dots, n$, taken from the alphabet \mathbb{Z}^d . If \mathcal{W}_n , $n = 1, 2, \dots$, denotes the set of words with n letters, then (10) may be rewritten as:

$$\sum_{w \in \mathcal{W}_n} \alpha_w(t) f_w(y_0). \quad (13)$$

If $\mathbf{k} \in \mathbb{Z}^d$ and $n = 1, 2, \dots$, the notation \mathbf{k}^n means $\mathbf{k} \cdots \mathbf{k} \in \mathcal{W}_n$. Two words $w = \mathbf{k}_1 \cdots \mathbf{k}_n$ and $w' = \mathbf{k}'_1 \cdots \mathbf{k}'_m$ may be concatenated to give rise to a new word $ww' = \mathbf{k}_1 \cdots \mathbf{k}_n \mathbf{k}'_1 \cdots \mathbf{k}'_m \in \mathcal{W}_{n+m}$. It is also convenient to introduce an empty word \emptyset such that $\emptyset w = w \emptyset = w$ for each w . The set of all words (including the empty word) will be denoted by \mathcal{W} .

With this notation, the expansion of y in (8) is, for each fixed t , an instance of the general format

$$W(\delta, y) = \delta_\emptyset y + \sum_{n=1}^{\infty} \epsilon^n \sum_{w \in \mathcal{W}_n} \delta_w f_w(y), \quad (14)$$

where δ is a map that associates with each $w \in \mathcal{W}$ a complex number δ_w . Series of the form (14) will be referred to as *word-series* (relative to equation (1)). In the particular case where the map δ is such that $\delta_\emptyset = 1$ and $\delta_w = 0$ for each $w \neq \emptyset$, the word-series is the identity $W(\delta, y) \equiv y$.

The sums (10) or (13) are not the only way of writing the coefficient of ϵ^n in the expansion of y . In [8, Sec. 2.2], y is expanded by reformulating the initial value problem (1)–(2) as an integral equation and then using Picard iteration. With that alternative approach, which does not use the Lie operators $E_{\mathbf{k}}$, in lieu of (13), one obtains a sum of the form

$$\sum_{u \in \mathcal{T}_n} \frac{1}{\sigma_u} \alpha_u(t) \mathcal{F}_u(y_0), \quad (15)$$

where the set of indices \mathcal{T}_n comprises the trees with n vertices mentioned above, and for each $u \in \mathcal{T}_n$, the integer σ_u is a normalizing factor, \mathcal{F}_u the corresponding elementary differential and $\alpha_u(t)$ a suitable complex coefficient, independent of f . The expressions (13) and (15) for the coefficient of ϵ^n in the expansion of y share of course a common value. Given (13), one obtains (15) by writing each word-basis function f_w as a linear combination of elementary differentials \mathcal{F}_u as in (11)–(12). Conversely, it is proved in [8, Sec. 3] (cf. [15]) that, given (15), one may reach (13) by grouping together different elementary differentials to form word-basis functions f_w .

Formal series in powers of ϵ

$$B(\eta, y) = \eta_\emptyset y + \sum_{n=1}^{\infty} \epsilon^n \sum_{u \in \mathcal{T}_n} \frac{1}{\sigma_u} \eta_u \mathcal{F}_u(y),$$

are called *B-series* and were first introduced in the analysis of Runge-Kutta numerical integration methods [12], [19], [11]. As we have just explained in the particular case of the word-series expansion of $y(t)$, each word-series may be rearranged to yield a B-series after

writing each word-basis function in terms of elementary differentials. However it is *not* true that each B-series may be rearranged to yield a word-series because elementary differentials cannot in general be expressed as linear combinations of word-basis functions;³ in particular there are B-series that are useful for averaging purposes and cannot be reformulated as word-series (see an instance in [8, Remark 3.3]). In this paper we shall emphasize the more compact word-series format that facilitates the derivation of error bounds. The results in the remainder of this section are all taken from [8], where they are formulated in the (more general) language of B-series.

2.3 The transport equation

By computing the iterated integrals in (7), we find, in particular,

$$\begin{aligned} \alpha_{\mathbf{k}}(t) &= \frac{i}{\mathbf{k} \cdot \boldsymbol{\omega}} (1 - e^{i\mathbf{k} \cdot \boldsymbol{\omega} t}), & \mathbf{k} \in \mathbb{Z} \setminus \{\mathbf{0}\}, \\ \alpha_{\mathbf{0}^n}(t) &= t^n / n!, \\ \alpha_{\mathbf{1}\mathbf{k}}(t) &= \frac{it}{\mathbf{1} \cdot \boldsymbol{\omega}} + \frac{1 - e^{i\mathbf{1} \cdot \boldsymbol{\omega} t}}{(\mathbf{1} \cdot \boldsymbol{\omega})^2}, & \mathbf{k} = -\mathbf{1} \neq \mathbf{0}, \\ \dots\dots &= \dots\dots\dots \end{aligned}$$

where we note that the left-hand sides are complex polynomials in the $2d + 1$ variables $t, e^{i\omega_1 t}, \dots, e^{i\omega_d t}, e^{-i\omega_1 t}, \dots, e^{-i\omega_d t}$. In general, it is easily proved by induction, that for each nonempty word w , there exists a unique complex function $\gamma(t, \theta)$, $t \in \mathbb{R}$, $\theta \in \mathbb{T}^d$, that is a polynomial in the variables $t, e^{i\theta_1}, \dots, e^{i\theta_d}, e^{-i\theta_1}, \dots, e^{-i\theta_d}$ and for which

$$\alpha_w(t) = \gamma_w(t, t\boldsymbol{\omega}). \tag{16}$$

This relation is extended to the empty word by defining $\alpha_{\emptyset}(t) \equiv 1$ and $\gamma_{\emptyset}(t, \theta) \equiv 1$.

Furthermore we have the following characterization [8, Sec. 2.4]:

Proposition 2.1 *For each $w \in \mathcal{W}$ and $\mathbf{k} \in \mathbb{Z}^d$, $\gamma_{w\mathbf{k}}$ is the only function that simultaneously satisfies the following requirements:*

- *It is a polynomial in the variables $t, e^{i\theta_1}, \dots, e^{i\theta_d}, e^{-i\theta_1}, \dots, e^{-i\theta_d}$.*
- *It is a solution of the problem*

$$(\partial_t + \boldsymbol{\omega} \cdot \nabla_{\boldsymbol{\theta}})\gamma_{w\mathbf{k}}(t, \boldsymbol{\theta}) = \gamma_w(t, \boldsymbol{\theta})e^{i\mathbf{k} \cdot \boldsymbol{\theta}}, \quad \gamma_{w\mathbf{k}}(0, \mathbf{0}) = 0. \tag{17}$$

The transport partial differential equation (17) plays in the formal series approach to averaging the role played by the homological equation in approaches via changes of variables. By using Proposition 2.1, one may derive the following recursions [8, Prop. 4.1] that allow the effective computation of the γ_w .

³It is possible to characterize the set of B-series that may be rearranged as word-series, see [8, Sec. 3], and [15]

Proposition 2.2 *If $n = 1, 2, \dots$, $\mathbf{k} \in \mathbb{Z}^d \setminus \{\mathbf{0}\}$ and $\mathbf{l}_1, \dots, \mathbf{l}_s \in \mathbb{Z}^d$, then:*

$$\begin{aligned}\gamma_{\mathbf{k}}(t, \theta) &= \frac{i}{\mathbf{k} \cdot \boldsymbol{\omega}} (1 - e^{i\mathbf{k} \cdot \theta}), \\ \gamma_{\mathbf{0}^n}(t, \theta) &= t^n / n!, \\ \gamma_{\mathbf{0}^n \mathbf{k}}(t, \theta) &= \frac{i}{\mathbf{k} \cdot \boldsymbol{\omega}} (\gamma_{\mathbf{0}^{n-1} \mathbf{k}}(t, \theta) - \gamma_{\mathbf{0}^n}(t, \theta) e^{i\mathbf{k} \cdot \theta}), \\ \gamma_{\mathbf{k} \mathbf{l}_1 \dots \mathbf{l}_s}(t, \theta) &= \frac{i}{\mathbf{k} \cdot \boldsymbol{\omega}} (\gamma_{\mathbf{l}_1 \dots \mathbf{l}_s}(t, \theta) - \gamma_{(\mathbf{k} + \mathbf{l}_1) \mathbf{l}_2 \dots \mathbf{l}_s}(t, \theta)), \\ \gamma_{\mathbf{0}^n \mathbf{k} \mathbf{l}_1 \dots \mathbf{l}_s}(t, \theta) &= \frac{i}{\mathbf{k} \cdot \boldsymbol{\omega}} (\gamma_{\mathbf{0}^{n-1} \mathbf{k} \mathbf{l}_1 \dots \mathbf{l}_s}(t, \theta) - \gamma_{\mathbf{0}^n (\mathbf{k} + \mathbf{l}_1) \mathbf{l}_2 \dots \mathbf{l}_s}(t, \theta)).\end{aligned}$$

It is also the transport equation that leads to the following result [8, Prop. 2.9], which is the key to later developments:

Proposition 2.3 *For each $t, t' \in \mathbb{R}$ and $y \in \mathbb{R}^D$*

$$W(\gamma(t', \mathbf{0}), W(\gamma(t, \mathbf{0}), y)) = W(\gamma(t + t', \mathbf{0}), y) \quad (18)$$

and for each $t \in \mathbb{R}$, $\theta \in \mathbb{T}^d$ and $y \in \mathbb{R}^D$

$$W(\gamma(0, \theta), W(\gamma(t, \mathbf{0}), y)) = W(\gamma(t, \theta), y) \quad (19)$$

2.4 Quasi-stroboscopic averaging

For $w \in \mathcal{W}$, let us define coefficient

$$\bar{\alpha}_w(t) := \gamma_w(t, 0) \quad (20)$$

and consider the family of transformations, parameterized by t ,

$$y \in \mathbb{R}^D \mapsto W(\bar{\alpha}(t), y) \in \mathbb{R}^D.$$

For $t = 0$, by (20) and (16), $W(\bar{\alpha}(0), y) = W(\gamma(0, \mathbf{0}), y) = W(\alpha(0), y) = y$, and the corresponding transformation is the identity. Then (18) shows that the family is a group and, by implication, the solution flow of an *autonomous* differential system (the averaged system). In order to write down such an averaged system, we differentiate $W(\bar{\alpha}(t), y)$ with respect to t at $t = 0$ and obtain

$$\frac{d}{dt} Y = W(\bar{\beta}, Y)$$

where the coefficients are given by

$$\bar{\beta}_w = \left. \frac{d}{dt} \bar{\alpha}_w(t) \right|_{t=0}, \quad w \in \mathcal{W}. \quad (21)$$

After recalling the definition of word-series in (14), we have, more explicitly,

$$\frac{d}{dt} Y = \epsilon F(Y, \epsilon), \quad F(Y) = F_1(Y) + \epsilon F_2(Y) + \dots + \epsilon^{n-1} F_n(Y) + \dots, \quad (22)$$

with

$$F_n(y) := \sum_{w \in \mathcal{W}_n} \bar{\beta}_w f_w(Y) \quad n = 1, 2, \dots \quad (23)$$

In particular, from Proposition 2.2, $\bar{\alpha}_{\mathbf{k}}(t) = 0$ for $\mathbf{k} \neq \mathbf{0}$ and $\bar{\alpha}_{\mathbf{0}} = t$, so that, from (21), $\bar{\beta}_{\mathbf{k}}(t) = 0$ for $\mathbf{k} \neq \mathbf{0}$ and $\bar{\beta}_{\mathbf{0}} = 1$ and therefore

$$F_1(Y) = f_{\mathbf{0}}(Y); \quad (24)$$

thus F_1 is the average of $f(\cdot, \theta)$ over $\theta \in \mathbb{T}^d$.

We now exploit the identity (19). If we define coefficients

$$\kappa_w(\theta) := \gamma_w(0, \theta), \quad w \in \mathcal{W}, \theta \in \mathbb{T}^d, \quad (25)$$

then (19) and (20) imply

$$W(\kappa(\theta), W(\bar{\alpha}(t), y)) = W(\gamma(t, \theta), y),$$

and in particular, for $\theta = t\omega$,

$$W(\kappa(t\omega), W(\bar{\alpha}(t), y)) = W(\gamma(t, t\omega), y) = W(\alpha(t), y).$$

Since $W(\alpha(t), y)$ and $W(\bar{\alpha}(t), y)$ are respectively the solution operator of the given oscillatory initial value problem (1)–(2) and the solution flow of the autonomous (22), we conclude that the change of variables $y(t) = W(\kappa(t\omega), Y(t))$ maps solutions of (22) onto solutions of (1). Note that $y(0) = Y(0)$ since, from (25), $\kappa_w(\mathbf{0}) = \gamma_w(0, \mathbf{0}) = 0$ for $w \neq \emptyset$ and $\kappa_{\emptyset}(\mathbf{0}) = \gamma_{\emptyset}(0, \mathbf{0}) = 1$. In this way we have proved the following result.

Theorem 2.4 *The solution of (1)–(2) may be written as*

$$y(t) = U(Y(t), t\omega, \epsilon),$$

where U is the change of variables parameterized by $\theta \in \mathbb{T}^d$

$$y = Y + \epsilon \check{U}(Y, \theta, \epsilon); \quad \check{U}(Y, \theta, \epsilon) := u_1(Y, \theta) + \dots + \epsilon^{n-1} u_n(Y, \theta) + \dots \quad (26)$$

with

$$u_n(Y, \theta) = \sum_{w \in \mathcal{W}_n} \kappa_w(\theta) f_w(Y), \quad n = 1, 2, \dots \quad (27)$$

and $Y(t)$ is the solution of the autonomous (averaged) system (22) with initial condition $Y(0) = y_0$.

We emphasize that the averaged system (22)–(23) is made up of two components: the coefficients $\bar{\beta}_w$ that are independent of f and the word-basis functions f_w . In fact the values of $\bar{\beta}_w$ can easily be found recursively, as shown in our next result, that is proved by differentiating with respect to t the formulae in Proposition 2.2 (see (20)–(21)).

Proposition 2.5 Given $n = 1, 2, \dots$, $\mathbf{k} \in \mathbb{Z}^d \setminus \{0\}$, and $\mathbf{l}_1, \dots, \mathbf{l}_s \in \mathbb{Z}^d$,

$$\begin{aligned}\bar{\beta}_{\mathbf{k}} &= 0, \\ \bar{\beta}_{\mathbf{0}} &= 1, \\ \bar{\beta}_{\mathbf{0}^{n+1}} &= 0, \\ \bar{\beta}_{\mathbf{0}^n \mathbf{k}} &= \frac{i}{\mathbf{k} \cdot \boldsymbol{\omega}} (\bar{\beta}_{\mathbf{0}^{n-1} \mathbf{k}} - \bar{\beta}_{\mathbf{0}^n}), \\ \bar{\beta}_{\mathbf{k} \mathbf{l}_1 \dots \mathbf{l}_s} &= \frac{i}{\mathbf{k} \cdot \boldsymbol{\omega}} (\bar{\beta}_{\mathbf{l}_1 \dots \mathbf{l}_s} - \bar{\beta}_{(\mathbf{k} + \mathbf{l}_1) \mathbf{l}_2 \dots \mathbf{l}_s}), \\ \bar{\beta}_{\mathbf{0}^n \mathbf{k} \mathbf{l}_1 \dots \mathbf{l}_s} &= \frac{i}{\mathbf{k} \cdot \boldsymbol{\omega}} (\bar{\beta}_{\mathbf{0}^{n-1} \mathbf{k} \mathbf{l}_1 \dots \mathbf{l}_s} - \bar{\beta}_{\mathbf{0}^n (\mathbf{k} + \mathbf{l}_1) \mathbf{l}_2 \dots \mathbf{l}_s}).\end{aligned}$$

The values of the coefficients κ_w in (25) that are required to write down the change of variables (26)–(27) may be recursively computed from Proposition 2.2. In fact in the present approach, the averaged system and the change of variables may be computed independently of each other, something that does not happen in standard approaches that use successive changes of variables.

In the periodic case ($d = 1$) the oscillatory and averaged solutions coincide at times $t = 2\pi j/\omega$, j integer, because they coincide at $t = 0$ and the change of variables $y = U(Y, \theta, \epsilon)$ is 2π periodic in θ . Therefore (22)–(23) provides the so-called *stroboscopic* averaged system for (1). Alternative averaged systems are discussed in detail in [8, Sec. 2.6]. In the quasi-periodic $d > 1$ case, the stroboscopic effect is not present: $y(t)$ and $Y(t)$ only coincide at time $t = 0$, because, due to non-resonance, the mapping $t \mapsto t\boldsymbol{\omega} \in \mathbb{T}^d$ only visits $\mathbf{0} \in \mathbb{T}^d$ for $t = 0$. In this situation we say that (22)–(23) is *quasi-stroboscopic* averaged system for (1).

To conclude this subsection let us discuss the assumptions on f . The hypothesis that the Fourier expansion consists of a finite number of (nontrivial) terms has been used to ensure that for each $n = 1, 2, \dots$ the series (10), (23), (27) converge and that, accordingly, (8), (22), (26) are bona fide formal series in powers of the parameter ϵ . Then the expansion (8) and Theorem 2.4 have to be understood in the sense of formal series. By truncating the formal series one of course obtains the corresponding Taylor polynomials. If those polynomials are only required up to a target maximum degree, then the differentiability requirements on f may be decreased correspondingly. It is also possible to work with f defined in $\mathcal{K} \times \mathbb{T}^d$, with \mathcal{K} a domain of \mathbb{R}^D , rather than with f defined on the whole of $\mathbb{R}^D \times \mathbb{T}^d$.

On the other hand the convergence of the series in (10), (23), (27) may be guaranteed with hypotheses on f other than the requirement that f has finitely-many non-vanishing Fourier coefficients, see e.g. Theorems 3.1 and 3.2 below.

2.5 An example

The system ($D = 2$, $d = 1$, $\omega = 1$)

$$\frac{d}{dt} y^1 = \epsilon, \quad \frac{d}{dt} y^2 = \epsilon \cos t g(y^1), \quad (28)$$

has three nontrivial Fourier coefficients

$$f_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad f_{+1} = f_{-1} = \frac{1}{2} \begin{bmatrix} 0 \\ g(y^1) \end{bmatrix}.$$

From (9), it follows that, for words $w \in W_{n+1}$, $n = 1, 2, \dots$, the basis function f_w vanishes identically except in cases where $w = 0^n k$, $k = \pm 1$. For these words $f_{0^n k} = (1/2)[0, g^{(n)}(y^1)]^T$. Furthermore, from Proposition 2.5, $\bar{\beta}_{0^n k} = -(i/k)^n$, $k = \pm 1$, and then the coefficients F_n in (23) are given by

$$F_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad F_{2j} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad F_{2j+1} = (-1)^{j+1} \begin{bmatrix} 0 \\ g^{(2j)}(y^1) \end{bmatrix}, \quad j = 1, 2, \dots$$

These lead to the following expression for the averaged system (22):

$$\frac{d}{dt}Y^1 = \epsilon, \quad \frac{d}{dt}Y^2 = \epsilon^3 g^{(2)}(Y^1) - \epsilon^5 g^{(4)}(Y^1) + \epsilon^7 g^{(6)}(Y^1) - \dots \quad (29)$$

The change of variables may be computed in a similar way.

2.6 Geometric properties

As proved in [8, Th. 3.2], for $n = 2, 3, \dots$, the terms in the series (23) for $F_n(y)$ can be rearranged to yield

$$F_n = \sum_{\mathbf{k}_1 \dots \mathbf{k}_n \in \mathcal{W}_n} \frac{1}{j} \beta_{\mathbf{k}_1 \dots \mathbf{k}_n} [[\dots [[f_{\mathbf{k}_1}, f_{\mathbf{k}_2}], f_{\mathbf{k}_3}] \dots], f_{\mathbf{k}_n}](y),$$

where, for each pair of vector fields, f, g ,

$$[f, g](y) = \left(\frac{\partial}{\partial y} g(y) \right) f(y) - \left(\frac{\partial}{\partial y} f(y) \right) g(y)$$

denotes the corresponding commutator (Lie bracket). This fact, in tandem with (24), implies that if the Fourier coefficients $f_{\mathbf{k}}$, $\mathbf{k} \in \mathbb{Z}^d$, belong to a specific Lie subalgebra of the Lie algebra of vector fields (e.g. each $f_{\mathbf{k}}$ is Hamiltonian) then the quasi-stroboscopic averaged (22) system will also belong to the same Lie subalgebra (e.g. averaging a Hamiltonian system will lead to a Hamiltonian system). Similarly the change of variables (26) will belong to the corresponding Lie group (in the example, the change of variables will be canonical).

Remark 2.6 Similarly the fact that the F_n can be written in terms of Lie brackets shows that if we work with a differential equation defined on a manifold (rather than in \mathbb{R}^D or in a domain $\mathcal{K} \subset \mathbb{R}^D$) the averaged equation will also be (globally) defined on that manifold.

By computing explicitly the coefficients by means of Proposition 2.5 and writing the F_n in terms of commutators, we find the following expression for the lower-order terms of the quasi-stroboscopically averaged system

$$\frac{d}{dt}Y = \epsilon F_1(Y) + \epsilon^2 F_2(Y) + \epsilon^3 [F_3^+(Y) + F_3^-(Y)] + \mathcal{O}(\epsilon^4),$$

where $F_1 = f_0$, and

$$\begin{aligned}
F_2 &= \sum_{\mathbf{k} > \mathbf{0}} \frac{i}{\mathbf{k} \cdot \omega} ([f_{-\mathbf{k}}, f_{\mathbf{k}}] + [f_{\mathbf{k}} - f_{-\mathbf{k}}, f_0]) \\
F_3^+ &= \sum_{\mathbf{k} > \mathbf{0}} \frac{1}{(\mathbf{k} \cdot \omega)^2} \left([f_0, [f_0, f_{\mathbf{k}}]] + [f_{\mathbf{k}}, [f_{\mathbf{k}}, f_{-\mathbf{k}}]] + \frac{1}{2} [f_{\mathbf{k}}, [f_0, f_{\mathbf{k}}]] + 2[f_{-\mathbf{k}}, [f_{\mathbf{k}}, f_0]] \right) \\
&\quad + \sum_{\mathbf{k} > \mathbf{1} > \mathbf{0}} \frac{1}{(\mathbf{k} \cdot \omega)(\mathbf{1} \cdot \omega)} \left([f_{-\mathbf{1}}, [f_{\mathbf{k}}, f_{\mathbf{1}-\mathbf{k}}]] - [f_{\mathbf{1}}, [f_{\mathbf{k}}, f_{-\mathbf{k}-\mathbf{1}}]] \right),
\end{aligned}$$

and F_3^- is obtained from F_3^+ by replacing $(\mathbf{k}, \mathbf{1})$ by $(-\mathbf{k}, -\mathbf{1})$. In these formulae, $<$ is a total ordering in the set of multi-indices \mathbb{Z}^d with the following property: the relations $\mathbf{k} > \mathbf{0}$ and $\mathbf{1} > \mathbf{0}$ imply $\mathbf{k} + \mathbf{1} > \mathbf{0}$.

3 Estimates

Our purpose is now to show how the present approach to averaging may be used to derive error bounds for the averaged system.

3.1 Preliminaries

In this section we assume that we are interested in studying the differential system (1) in a domain $\mathcal{K} \subset \mathbb{R}^D$. The symbol $\|\cdot\|$ will refer to a norm in \mathbb{C}^D or to the associated norm for $D \times D$ complex matrices. For $\rho \geq 0$ we denote

$$\mathcal{K}_\rho = \{y + z \in \mathbb{C}^D : y \in \bar{\mathcal{K}}, \|z\| \leq \rho\}$$

($\bar{\mathcal{K}}$ is the closure of \mathcal{K}) and for vector or matrix-valued bounded functions ϕ defined in \mathcal{K}_ρ we set

$$\|\phi\|_\rho = \sup_{y \in \mathcal{K}_\rho} \|\phi(y)\|.$$

Our hypotheses on f are now as follows:

Assumption 1 *There exist $R > 0$ and an open set $\mathcal{U} \supset \mathcal{K}_R$, such that, for each $\theta \in \mathbb{T}^d$, $f(\cdot, \theta)$ may be extended to a map $\mathcal{U} \rightarrow \mathbb{C}^D$ that is analytic at each point $y \in \mathcal{K}_R$. Furthermore the Fourier coefficients $f_{\mathbf{k}}$ of f have bounds*

$$\forall \mathbf{k} \in \mathbb{Z}^d, \quad \|f_{\mathbf{k}}\|_R \leq a_{\mathbf{k}}, \quad a_{\mathbf{k}} \geq 0,$$

with

$$M := \sum_{\mathbf{k} \in \mathbb{Z}^d} a_{\mathbf{k}} < \infty.$$

Under this assumption, the Fourier series (6) converges absolutely and uniformly in $\mathcal{K}_R \times \mathbb{T}^d$ and therefore f is (jointly) continuous in $\mathcal{K}_R \times \mathbb{T}^d$. Furthermore

$$\forall \theta \in \mathbb{T}^d, \quad \|f(\cdot, \theta)\|_R \leq M.$$

On the other hand the assumption is not strong enough to guarantee the differentiability of f with respect to θ .

3.2 Convergence of the expansion of y

Our first result, while not being necessary to prove the estimates in Theorem 3.4, provides an example of the techniques employed in this section:

Theorem 3.1 *Suppose that f satisfies the requirements in Assumption 1.*

1. *For $0 \leq \rho < R$, $n = 1, 2, \dots$, $y_0 \in \mathcal{K}_\rho$ and $t \in \mathbb{R}$, the series (10) converges absolutely. The convergence is uniform in y_0 .*
2. *For $0 \leq \rho < R$, $y_0 \in \mathcal{K}_\rho$ and $|\epsilon t| < (R - \rho)/(eM)$, the expansion of $y(t)$ in powers of ϵ in (8) is absolutely convergent. The convergence is uniform for $y_0 \in \mathcal{K}_\rho$ and $|\epsilon t|$ ranging in any compact subinterval of the interval $[0, (R - \rho)/(eM)]$.*

Proof: Propositions 4.1 and 4.5 show that for, $n \geq 2$,

$$\begin{aligned} \sum_{w \in \mathcal{W}_n} \|\epsilon^n \alpha_w(t) f_w(y_0)\| &\leq \frac{|\epsilon t|^n (n-1)^{n-1}}{n! (R-\rho)^{n-1}} \sum_{\mathbf{k}_1, \dots, \mathbf{k}_n \in \mathbb{Z}^d} a_{\mathbf{k}_1} \cdots a_{\mathbf{k}_n} \\ &= \frac{|\epsilon t|^n (n-1)^{n-1}}{n! (R-\rho)^{n-1}} M^n. \end{aligned}$$

The series whose n -th term is the last expression converges for $|\epsilon t| < (R - \rho)/(eM)$. \square

It may be pointed out that, since $\|f\| \leq M$ in $\mathcal{K}_R \times \mathbb{T}^d$, the solution of (1)–(2) with $y_0 \in \mathcal{K}_\rho$ exists at least in an interval of length $(R - \rho)/(M\epsilon)$. This length is of course larger than the time-span $(R - \rho)/(eM\epsilon)$ that features in the theorem.

3.3 Exponentially small estimates

We now address the derivation of error estimates. Only periodic case $d = 1$ will be considered in the estimates of the present article. Note that in the periodic situation we may always rescale t so as to have $\omega = 1$. The next result refers to the coefficients F_n and u_n in (23) and (27).

Theorem 3.2 *Suppose that f satisfies the requirements in Assumption 1, $d = 1$, $\omega = 1$. For $n = 1, 2, \dots$, $0 \leq \rho < R$, $y \in \mathcal{K}_\rho$, $\theta \in \mathbb{T}$, the series in (23) and (27) are absolutely and uniformly convergent. Furthermore, the functions $F_n(y)$, $u_n(y, \theta)$ defined by those series satisfy⁴*

$$\|F_n\|_\rho \leq \frac{1}{2} \frac{(2M)^n (n-1)^{n-1}}{(R-\rho)^{n-1}}, \quad (30)$$

$$\|u_n(\cdot, \theta)\|_\rho \leq \frac{(2M)^n (n-1)^{n-1}}{(R-\rho)^{n-1}}, \quad (31)$$

$$\|\partial_\theta u_n(\cdot, \theta)\|_\rho \leq \frac{1}{2} \frac{(2M)^n (n-1)^{n-1}}{(R-\rho)^{n-1}}, \quad (32)$$

$$\|\partial_y u_n(\cdot, \theta)\|_\rho \leq \frac{(2M)^n n^n}{(R-\rho)^n}. \quad (33)$$

⁴Throughout the paper it is understood that for $n = 1$ the expression $(n - 1)^{n-1}$ takes the value 1.

Proof: It is based on Propositions 4.1 and 4.5 and very similar to that of Theorem 3.1; details will not be given. Note that $\partial_\theta u_n$ and $\partial_y u_n$ are represented by the series obtained by differentiating term by term the series for u_n . \square

Since the series with n -th term

$$\epsilon^n \frac{(2M)^n (n-1)^{n-1}}{(R-\rho)^{n-1}}$$

diverges for $\epsilon \neq 0$, the bounds (30) and (31) do not allow us to prove the convergence of the series for the averaged system and change of variables in (22) and (26). In fact it is well known that these series are in general divergent⁵ and the time dependence of (1) cannot be completely eliminated by means of a change of variables. In our set-up it is easy to construct an explicit example of divergence. Consider the system (28) with

$$g(y^1) = \frac{1}{1-y^1};$$

Assumption 1 holds on any domain $\mathcal{K} \subset \mathbb{R}^2$ whose closure $\bar{\mathcal{K}}$ does not intersect the line $\{y^1 = 1\} \subset \mathbb{R}^2$. The series in (29) has $g^{(n)}(y^1) = n!/(1-y^1)^{n+1}$ and it is therefore divergent for each $\epsilon \neq 0$, $y \in \mathcal{K}$; the series for the change of variables is found to diverge similarly.

The divergence of the series (22) and (26) entails that they have to be truncated in order to describe approximately the dynamics of (1). Our next result presents some properties of truncated changes of variables. Item 1. implies that the change of variables is an $\mathcal{O}(\epsilon)$ perturbation of the identity. Item 3. guarantees the local invertibility of the change. Item 4. shows global invertibility under the assumption of convexity. This assumption is not necessary, as global invertibility could be shown after finding the corresponding inverse transformation, something that may be achieved by solving for Y equation (34) with the help of fixed-point iteration.

Theorem 3.3 *Suppose that f satisfies the requirements in Assumption 1, $d = 1$, $\omega = 1$. For $N = 1, 2, \dots$ consider the change of variables*

$$y = Y + \epsilon \check{U}^{(N)}(Y, t, \epsilon) \quad (34)$$

with

$$\check{U}^{(N)}(y, \theta, \epsilon) := u_1(y, \theta) + \epsilon u_2(y, \theta) + \dots + \epsilon^{N-1} u_N(y, \theta)$$

(the functions u_n are as in the preceding theorem). Assume that $\epsilon \in \mathbb{C}$ satisfies:

$$|\epsilon| \leq \epsilon_0, \quad \epsilon_0 = \epsilon_0(N) := \frac{R}{8M} \frac{1}{N}, \quad (35)$$

then:

1. For each $\theta \in \mathbb{T}$, $\|\check{U}^{(N)}(\cdot, \theta, \epsilon)\|_{R/2} \leq 3M$ and $\|\partial_\theta \check{U}^{(N)}(\cdot, \theta, \epsilon)\|_{R/2} \leq 3M/2$.
2. For each $\theta \in \mathbb{T}$, the mapping $Y \in \mathcal{K}_{R/2} \mapsto Y + \epsilon \check{U}^{(N)}(Y, \theta, \epsilon)$ is analytic and takes values in \mathcal{K}_R .

⁵See however Section 3.4 below.

3. For each $\theta \in \mathbb{T}$ and $Y \in \mathcal{K}_{R/2}$, the Jacobian matrix $I + \epsilon \partial_Y \check{U}^{(N)}(Y, \theta, \epsilon)$ is invertible and

$$\|(I + \epsilon \partial_Y \check{U}^{(N)})^{-1}\| \leq 2.$$

4. If $\mathcal{K}_{R/2}$ is convex then, for each $\theta \in \mathbb{T}$, the mapping $Y \in \mathcal{K}_{R/2} \mapsto Y + \epsilon \check{U}^{(N)}(Y, \theta, \epsilon)$ is one-to-one.

Proof: From (31) with $\rho = R/2$, (35) and Lemma 4.7:

$$\begin{aligned} \|\check{U}^{(N)}(\cdot, t, \epsilon)\|_{R/2} &\leq 2M \sum_{n=1}^N |\epsilon|^{n-1} \frac{(4M)^{n-1} (n-1)^{n-1}}{R^{n-1}} \\ &\leq 2M \sum_{n=1}^N \frac{(n-1)^{n-1}}{(2N)^{n-1}} \leq 3M. \end{aligned}$$

Then, for $Y \in \mathcal{K}_{R/2}$, $\|\epsilon \check{U}^{(N)}(Y, t, \epsilon)\| \leq 3RM/(8MN) < R/2$ and, accordingly, $y \in \mathcal{K}_R$.

The derivation of the bound for $\partial_\theta \check{U}^{(N)}$, based on (32), is similar to the derivation above of the estimate for $\check{U}^{(N)}$.

From (33) with $\rho = R/2$, (35) and Lemma 4.7, for $Y \in \mathcal{K}_{R/2}$,

$$\|\epsilon \partial_Y \check{U}^{(N)}(Y, t, \epsilon)\| \leq \sum_{n=1}^N |\epsilon|^n \frac{(4M)^n n^n}{R^n} \leq \sum_{n=1}^N \frac{n^n}{(2N)^n} \leq \frac{1}{2}.$$

Standard results show that $I + \epsilon \partial_Y \check{U}^{(N)}(Y, t, \epsilon)$ is then invertible and the norm of its inverse is $\leq 1/(1 - 1/2) = 2$.

If $\mathcal{K}_{R/2}$ is convex and $Y_1 + \epsilon \check{U}^{(N)}(Y_1, \theta, \epsilon) = Y_2 + \epsilon \check{U}^{(N)}(Y_2, \theta, \epsilon)$, then by the mean-value theorem

$$\|Y_1 - Y_2\| = \|\epsilon \check{U}^{(N)}(Y_1, \theta, \epsilon) - \epsilon \check{U}^{(N)}(Y_2, \theta, \epsilon)\| \leq \frac{1}{2} \|Y_1 - Y_2\|$$

and $Y_1 = Y_2$. \square

We are now in a position to establish our main result:

Theorem 3.4 *Suppose that f satisfies the requirements in Assumption 1, $d = 1$, $\omega = 1$. The application of the change of variables in Theorem 3.3 subject to (35) to the initial value problem (1)–(2) results in a problem*

$$\frac{d}{dt} Y = \epsilon (F^{(N)}(Y, \epsilon) + R^{(N)}(Y, t, \epsilon)), \quad Y(0) = y_0, \quad (36)$$

where

$$F^{(N)}(y, \epsilon) = F_1(y) + \epsilon F_2(y) + \dots + \epsilon^{N-1} F_N(y)$$

(the functions F_j are as defined in Theorem 3.2). The averaged vector field $F^{(N)}$ and the remainder $R^{(N)}$ possess the bounds

$$\|F^{(N)}(\cdot, \epsilon) - f_0(\cdot)\|_{R/2} \leq \frac{M}{2} \epsilon$$

and

$$\|R^{(N)}(\cdot, t, \epsilon)\|_{R/2} \leq \frac{5(\epsilon/\epsilon_0)^N}{1 - (\epsilon/\epsilon_0)} M. \quad (37)$$

In particular, assume that for given ϵ , with $|\epsilon| \leq R/(8eM)$, N is chosen as the integer part of the real number $R/(8eM\epsilon) \geq 1$. Then

$$\|R^{(N)}(Y, \theta, \epsilon)\|_{R/2} \leq \frac{5e^2}{e-1} M \exp\left(-\frac{R}{8eM} \frac{1}{\epsilon}\right).$$

Proof: The change of variables leads to the differential equation

$$\frac{d}{dt} Y = \epsilon g^{(N)}(Y, t, \epsilon),$$

where

$$g^{(N)}(Y, \theta, \epsilon) = (I + \epsilon \partial_Y \check{U}^{(N)}(Y, \theta, \epsilon))^{-1} (f(Y + \epsilon \check{U}^{(N)}(Y, \theta, \epsilon), \theta) - \partial_\theta \check{U}^{(N)}(Y, \theta, \epsilon)).$$

We then define $R^{(N)}(Y, \theta, \epsilon) = g^{(N)}(Y, \theta, \epsilon) - F^{(N)}(Y, \epsilon)$ with $F^{(N)}$ as in the statement of the theorem. Clearly, from the results on formal series in the preceding section, $F^{(N)}$ is the Taylor polynomial of degree $N - 1$ of $g^{(N)}$ seen as a function of ϵ . By standard Cauchy estimates in the complex disk of center $0 \in \mathbb{C}$ and radius ϵ_0 , we may write

$$\|R^{(N)}(Y, \theta, \epsilon)\|_{R/2} \leq \sum_{n=N}^{\infty} (\epsilon/\epsilon_0)^n \|g^{(N)}(\cdot, \theta, \epsilon_0)\|_{R/2} = \frac{(\epsilon/\epsilon_0)^N}{1 - (\epsilon/\epsilon_0)} \|g^{(N)}(\cdot, \theta, \epsilon_0)\|_{R/2}.$$

For $Y \in \mathcal{K}_{R/2}$, item 2. in Theorem 3.3 implies that $Y + \epsilon \check{U}^{(N)}(Y, t, \epsilon) \in \mathcal{K}_R$ and therefore

$$\|f(Y + \epsilon \check{U}^{(N)}(Y, t, \epsilon), t)\| \leq M.$$

Items 1. and 3. in the same Theorem then show that $\|g^{(N)}(\cdot, \theta, \epsilon_0)\|_{R/2} \leq 2(M + 3M/2)$. This establishes the bound for $R^{(N)}$ in (37). The bound for $F^{(N)}$ is derived as the bound for $\check{U}^{(N)}$ in Theorem 3.3.

With the choice of N given in the theorem, $\epsilon/\epsilon_0 \leq 1/e$ and therefore, from (37):

$$\|R^{(N)}(Y, \theta, \epsilon)\|_{R/2} \leq \frac{5}{1 - (1/e)} M \exp(-N).$$

Since $N > R/(8eM\epsilon) - 1$, the proof is complete. \square

Remark 3.5 Since f only enters the error bounds through the values of R and M , Theorem 3.4 also applies to the case where $f = f(y, \theta; \epsilon)$ depends on the parameter ϵ , provided that it satisfies Assumption 1 with R, M independent of ϵ .

Remark 3.6 In the geometric scenario, where f belongs to a specific Lie subalgebra of vector fields, the *truncated* transformed vector field $\epsilon F^{(N)}$ belongs to the same Lie subalgebra because, as we know, the F_n do. This is important in many applications, for instance in proving Proposition 2 in [16] on the preservation of the value of the Hamiltonian function (energy)

over exponentially long time intervals. On the other hand, since the change of variables in Theorem 3.3 has been chosen to be a polynomial in ϵ , it does not possess the favorable geometric properties of the formal, not truncated change discussed in the preceding section. For instance the truncated change will not be a canonical transformation if the given oscillatory system is Hamiltonian; by implication the vector field $\epsilon(F^{(N)} + R^{(N)})$ in (36) will not be Hamiltonian, as distinct from the situation with the truncated $\epsilon F^{(N)}$.⁶

3.4 The linear case

As discussed above, under Assumption 1, the series for the averaged system and change of variables in (22) and (26) are in general divergent. However, in particular cases, (22) and (26) converge; in those case the time-dependence of (1) may be eliminated completely by changing variables. As an example we treat the linear situation. We need the following hypothesis:

Assumption 2 *Assume that for each $\mathbf{k} \in \mathbb{Z}^d$, $f_{\mathbf{k}}(y) = A_{\mathbf{k}}y$, where the (constant) matrices $A_{\mathbf{k}}$ are such that, if $a_{\mathbf{k}} = \|A_{\mathbf{k}}\|$, then*

$$M := \sum_{\mathbf{k} \in \mathbb{Z}^d} a_{\mathbf{k}} < \infty.$$

This requirement implies that the Fourier series (6) converges absolutely and uniformly for y in an arbitrary ball in \mathbb{R}^D and $\theta \in \mathbb{T}^d$. Furthermore $\|f(y, \theta)\| \leq M\|y\|$ for $y \in \mathbb{R}^D$ and $\theta \in \mathbb{T}^d$.

Under Assumption 2 we may invoke the bounds for the word-basis functions in Proposition 4.6 instead of the weaker bounds in Proposition 4.5. Accordingly Theorem 3.2 may be strengthened to obtain:

Theorem 3.7 *Suppose that f satisfies the requirements in Assumption 2, $d = 1$, $\omega = 1$. For $n = 1, 2, \dots$, $y \in \mathbb{R}^D$, $\theta \in \mathbb{T}$, the series in (23) and (27) are absolutely convergent. The convergence is uniform for y in an arbitrary ball $\subset \mathbb{R}^D$ and $\theta \in \mathbb{T}$. Furthermore, the functions $F_n(y)$, $u_n(y, \theta)$ defined by those series satisfy*

$$\begin{aligned} \|F_n(y)\| &\leq \frac{1}{2}(2M)^n \|y\|, \\ \|u_n(y, \theta)\|_{\rho} &\leq (2M)^n \|y\|. \end{aligned}$$

The series in powers of ϵ for the averaged system and change of variables in (22) and (26) converge provided that $|\epsilon| < 1/(2M)$. The convergence is uniform for y in an arbitrary ball $\subset \mathbb{R}^D$ and $\theta \in \mathbb{T}$.

The convergence of (22) and (26) in the linear case has been proved in [6] (see also [3]) using the Floquet-Magnus expansion. It is remarkable that if one tries to bound the radius of convergence of (22) and (26) in terms of M by applying known results from the convergence of Magnus expansions [3] then one arrives at the condition $|\epsilon| < 1/(2M)$ that features in the preceding theorem.

⁶It is well known that it is possible to perturb the change in Theorem 3.3 with $\mathcal{O}(\epsilon^{N+1})$ terms in order to guarantee that the transformation $Y \mapsto y$ belongs to the relevant Lie subgroup.

Remark 3.8 The material in this subsection illustrates the way in which information on f may be used to derive information on the word-basis functions f_w that, in turn, can be incorporated directly into the analysis in view of the universal f -independent character of the coefficients $\bar{\beta}$ and κ .

4 Auxiliary results

This section gathers a number of technical details. Since the averaged vector field and the change of variables are built up in terms of the coefficients γ_w and the word-basis functions f_w , we need to estimate these two ingredients, something we do in Propositions 4.1 and 4.5 respectively.

Proposition 4.1 *The coefficients α_w defined in (7) satisfy $|\alpha_{\mathbf{k}_1 \dots \mathbf{k}_n}(t)| \leq t^n/n!$ for each $n = 1, 2, \dots, t \in \mathbb{R}$.*

Additionally, if $d = 1$ and $\omega = 1$, then:

1. *The coefficients κ_w in (25) satisfy $|\kappa_w(\theta)| \leq 2^n$ for $n = 1, 2, \dots, w \in \mathcal{W}_n, \theta \in \mathbb{T}$.*
2. *For $n = 1, 2, \dots, w \in \mathcal{W}_n$ and $\theta \in \mathbb{T}$, $|\partial_\theta \kappa(\theta)| \leq 2^{n-1}$.*
3. *The coefficients $\bar{\beta}_w$ in (21) satisfy $|\bar{\beta}_w| \leq 2^{n-1}$, for $n = 1, 2, \dots, w \in \mathcal{W}_n$.*

Proof: The bound for the coefficients α is trivial from the definition in (7). To bound $\kappa_w(\theta) = \gamma_w(0, \theta)$ and $\bar{\beta}_w$ use induction in the formulae in Proposition 2.2 and 2.5 respectively. For item 2, differentiate first with respect to θ the recursions in Proposition 2.2 to obtain recursions for the values $\partial_\theta \kappa(\theta) = \partial_\theta \gamma_w(0, \theta)$; then use induction. \square

Remark 4.2 Note that, in the quasi-periodic case, the recursions in Proposition 2.2 contain the small divisors; accordingly the task of bounding the coefficients is more delicate in that case than in the periodic scenario addressed here. In our approach, this is the only point where the treatments of the periodic and quasi-periodic cases differ from each other.

The proof of Proposition 4.5 will use the following two lemmas (cf. [11, Chap. IX, Lemma 7.4]).

Lemma 4.3 *For $\rho \geq 0, \delta > 0$, let g_1 (resp. g_2) be a \mathbb{C}^D -valued function analytic at each point in $\mathcal{K}_{\rho+\delta}$ (resp. \mathcal{K}_ρ). Then*

$$\|(\partial_y g_1)g_2\|_\rho \leq \frac{\|g_1\|_{\rho+\delta}\|g_2\|_\rho}{\delta}.$$

Proof: The chain rule implies that

$$\partial_y g_1(y)g_2(y) = \frac{d}{d\tau} g_1(y + \tau g_2(y)) \Big|_{\tau=0}.$$

Then the result is a direct consequence of the Cauchy bound for the derivative of the mapping $\tau \mapsto g_1(y + \tau g_2(y))$ that is analytic in the disk $|\tau| \leq \delta/\|g_2\|_\rho$. \square

Lemma 4.4 Let A be a function defined on \mathcal{K}_ρ , $\rho \geq 0$ with values in the set of $D \times D$ complex matrices. If $\|A(\cdot)v\|_\rho := \sup_{y \in \mathcal{K}_\rho} \|A(y)v\| \leq C$ for each vector $v \in \mathbb{C}^D$ with $\|v\| = 1$, then $\|A\|_\rho := \sup_{y \in \mathcal{K}_\rho} \|A(y)\| \leq C$.

Proof: To each $y \in \mathcal{K}_\rho$ there corresponds $v = v_y \in \mathbb{C}^D$ with $\|A(y)\| = \|A(y)v_y\|$, therefore, by hypothesis, $\|A(y)\| \leq C$. \square

Proposition 4.5 If f satisfies Assumption 1 then, for $0 \leq \rho < R$, $n = 1, 2, \dots$ and $\mathbf{k}_1, \dots, \mathbf{k}_n \in \mathbb{Z}^d$

$$\|f_{\mathbf{k}_1 \dots \mathbf{k}_n}\|_\rho \leq \frac{(n-1)^{n-1}}{(R-\rho)^{n-1}} a_{\mathbf{k}_1} \cdots a_{\mathbf{k}_n} \quad (38)$$

and

$$\|\partial_y f_{\mathbf{k}_1 \dots \mathbf{k}_n}\|_\rho \leq \frac{n^n}{(R-\rho)^n} a_{\mathbf{k}_1} \cdots a_{\mathbf{k}_n}. \quad (39)$$

Proof: From Lemma 4.3 with $\delta = (R-\rho)/n$ and the definition $f_{\mathbf{k}_1, \dots, \mathbf{k}_n}$ in (9) we may write, for $n = 1, 2, \dots$ and an arbitrary analytic map $g : \mathcal{K}_\rho \rightarrow \mathbb{C}^D$:

$$\begin{aligned} \|\partial_y f_{\mathbf{k}_1 \dots \mathbf{k}_n} g\|_\rho &\leq \frac{n}{R-\rho} \|f_{\mathbf{k}_1 \dots \mathbf{k}_n}\|_{\rho+(R-\rho)/n} \|g\|_\rho \\ &= \frac{n}{R-\rho} \|\partial_y f_{\mathbf{k}_2 \dots \mathbf{k}_n} f_{\mathbf{k}_1}\|_{\rho+(R-\rho)/n} \|g\|_\rho \\ &\leq \frac{n^2}{(R-\rho)^2} \|f_{\mathbf{k}_2 \dots \mathbf{k}_n}\|_{\rho+2(R-\rho)/n} a_{\mathbf{k}_1} \|g\|_\rho \\ &\leq \dots \\ &\leq \frac{n^n}{(R-\rho)^n} \|f_{\mathbf{k}_n}\|_{\rho+n(R-\rho)/n} a_{\mathbf{k}_1} \cdots a_{\mathbf{k}_{n-1}} \|g\|_\rho \\ &\leq \frac{n^n}{(R-\rho)^n} a_{\mathbf{k}_1} \cdots a_{\mathbf{k}_n} \|g\|_\rho. \end{aligned}$$

By choosing g to be a constant function and invoking Lemma 4.4 we obtain (39). Furthermore the choice $g = f_{\mathbf{k}_{n+1}}$ proves (38) with $n+1$ in lieu of n . This leaves us with the case $n = 1$ of (38) which is trivial from Assumption 1. \square

In the linear case the bounds may be improved. The proof of the following result is trivial:

Proposition 4.6 If f satisfies Assumption 2 then, for $n = 1, 2, \dots$ and $\mathbf{k}_1, \dots, \mathbf{k}_n \in \mathbb{Z}^d$, the constant matrix $\partial_y f_{\mathbf{k}_1 \dots \mathbf{k}_n} = A_{\mathbf{k}_n} \cdots A_{\mathbf{k}_1}$ satisfies:

$$\|\partial_y f_{\mathbf{k}_1 \dots \mathbf{k}_n}\| \leq a_{\mathbf{k}_1} \cdots a_{\mathbf{k}_n}$$

and therefore, for $y \in \mathbb{R}^D$,

$$\|f_{\mathbf{k}_1 \dots \mathbf{k}_n}(y)\| \leq a_{\mathbf{k}_1} \cdots a_{\mathbf{k}_n} \|y\|.$$

Finally, the following result was used in the proof of Theorem 3.3.

Lemma 4.7 For each $N = 1, 2, \dots$

$$b_N := \sum_{n=1}^N \frac{n^n}{(2N)^n} \leq \frac{1}{2}, \quad b_N^* := \sum_{n=1}^N \frac{(n-1)^{n-1}}{(2N)^{n-1}} \leq \frac{3}{2}.$$

Proof: The value $c_{n,N} = n^n/(2N)^n$ decreases as N increases. For $N \geq 5$ the decrease $c_{1,N} - c_{1,N+1} = 1/(2N(N+1))$ is larger than $c_{N+1,N+1}$ and therefore $b_{N+1} < b_N$. Since $b_1 = 1/2$, $b_2 = 1/2$, $b_3 = 29/72$, $b_4 = 155/512$, $b_5 = 4477/2000$, the sequence b_N is monotonically decreasing and bounded above by $1/2$. On the other hand $b_N^* \leq 1 + b_N$. \square

Acknowledgement. P. Chartier has received financial support from INRIA through the associated team MIMOL. A. Murua and J.M. Sanz-Serna have been supported by projects MTM2010-18246-C03-03 and MTM2010-18246-C03-01 respectively (Ministerio de Ciencia e Innovación).

References

- [1] V.I. Arnold, Geometrical Methods in the Theory of Ordinary Differential Equations, 2nd ed., Springer, New York, 1988.
- [2] V.I. Arnold, Mathematical Methods of Classical Mechanics, 2nd ed., Springer, New York, 1989.
- [3] S. Blanes, F. Casas, J.A. Oteo, J. Ros, The Magnus expansion and some of its applications, *Phys. Rep.* **470**, 151–238 (2009).
- [4] M.P. Calvo, Ph. Chartier, A. Murua, J.M. Sanz-Serna, A stroboscopic method for highly oscillatory problems. In: B. Engquist, O. Runborg, R. Tsai, (eds.) Numerical Analysis of Multiscale Computations, Lect. Notes Comput. Sci. Eng. 82, pp. 73–87. Springer, Berlin, 2011.
- [5] M.P. Calvo, Ph. Chartier, A. Murua, J.M. Sanz-Serna, Numerical stroboscopic averaging for ODEs and DAEs, Applied Numerical Mathematics, *Appl. Numer. Math.* **61**, 1077–1095 (2011).
- [6] F. Casas, J. A. Oteo, J. Ros, Floquet theory: exponential perturbative treatment, *J. Phys. A* **34**, 3379–3388 (2001).
- [7] P. Chartier, A. Murua, J.M. Sanz-Serna, Higher-Order averaging, formal series and numerical integration I: B-series, *Found. Comput. Math.* **10**, 695–727 (2010).
- [8] P. Chartier, A. Murua, J.M. Sanz-Serna, Higher-Order averaging, formal series and numerical integration II: the quasi-periodic case, *Found. Comput. Math.* submitted.
- [9] K.T. Chen, Integration of paths, geometric invariants and a generalized Baker-Hausdorff formula, *Annals Math.* **65**, 163–178 (1965).
- [10] M. Fliess, Fonctionelles causales nonlinéaires et indéterminées noncommutatives, *Bull. Soc. Math. France* **109**, 3–40 (1981).

- [11] E. Hairer, Ch. Lubich and G. Wanner, Geometric Numerical Integration, 2nd ed., Springer, Berlin, 2006.
- [12] E. Hairer, S.P. Nørsett and G. Wanner, Solving Ordinary Differential Equations I, Non-stiff Problems, 2nd ed., Springer, Berlin, 1993.
- [13] P. Lochak, C. Meunier, Multiphase Averaging for Classical Systems, with Applications to Adiabatic Theorems, Springer, New York, 1988.
- [14] A. Murua, Formal series and numerical integrators, part I: systems of ODEs and symplectic integrators, *Appl. Numer. Math.* **29**, 221–251 (1999).
- [15] A. Murua, The Hopf algebra of rooted trees, free Lie algebras and Lie series, *Found. Comput. Math.* **6**, 387–426 (2006).
- [16] A.I. Neishtadt, The separation of motions in systems with rapidly rotating phase, *J. Appl. Math. Mech.* **48**, 133–139 (1984).
- [17] J.-P. Ramis, R. Schäfke, Gevrey separation of fast and slow variables, *Nonlinearity* **9**, 353–384 (1996).
- [18] J.A. Sanders, F. Verhulst and J. Murdock, Averaging Methods in Nonlinear Dynamical Systems (2nd. ed.), Springer, New York, 2007.
- [19] J.M. Sanz-Serna and M. P. Calvo, Numerical Hamiltonian Problems, Chapman and Hall, London, 1994.
- [20] C. Simó, Averaging under fast quasi-periodic forcing. In: I. Seimenis (ed.) Proceedings of the NATO-ARW Integrable and Chaotic Behaviour in Hamiltonian Systems, Torun, Poland, 1993, pp. 13–34. Plenum, New York, 1994.
- [21] H. Sussman, A product expansion of the Chen series. In: C. Byrnes and A. Linnik (eds.) Applications of Nonlinear Control Systems, pp. 325–335. Elsevier, Amsterdam 1986.