# Reduced-Dimensional Models of Porous-Medium Convection

by

Navid Dianati Maleki

A dissertation submitted in partial fulfillment
of the requirements for the degree of
Doctor of Philosophy
(Physics)
in the University of Michigan
2013

Doctoral Committee:

Professor Charles R. Doering, Chair
Professor Mark E. J. Newman
Professor Leonard M. Sander
Professor Divakar Viswanath

# DEDICATION

To maman and baba.

# ACKNOWLEDGMENT

I would like to thank Charlie Doering for being an excellent expositor, an outstanding mentor, and a source of inspiration. I would also like to thank my collaborators Gregory Chini and Baole Wen for many fruitful discussions.

Finally, I am indebted to Elaine, without whose love and support, I would be lost.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF APPENDICES

# CHAPTER I

# Introduction

This dissertation summarizes the results of an investigation into low-dimensional modeling of an infinite-dimensional dynamical system, namely the problem of buoyancy-driven convection in a fluid-saturated porous medium. Motivated by the emergence of coherent structures in various dynamical regimes of this problem and the ultimate development of an "orderly" chaos, we ask whether it is possible to construct finite-dimensional dynamical models that are tailored specifically to this problem, reflecting its inherent symmetries and other qualitative features. We further ask if such models are in any sense more efficient in reproducing the essential features of the dynamics than other standard "generic" methods.

As we begin to study more and more complex dynamical regimes, we can not help but wonder if the coherent structures emerging amid the chaos can be exploited to further reduce the models in size. If successful, such an approach will constitute not only a computational advance, but a major step in identifying and isolating the "essence" of the motion in its most parsimonious form, thus bringing us one step closer to a physical understanding.

Thus, the effort is divided between two distinct but related fronts: on the one hand, we use direct numerical simulations in order to understand whether the most robust emergent coherent structures may in any way be seen as autonomous dynamical "units" encapsulating

the essential dynamics. On the other hand, we develop an *a priori* numerical method that yields a family of low-dimensional dynamical models of the problem which are adapted specifically to the equations of motion. Finally, we combine the findings of the two, and test the capabilities of the models enhanced by the physical insights gained from the study of coherent structures.

The rest of this dissertation is organized as follows: in Chapter II we present an overview of the theory of global attractors for driven dissipative dynamical systems, the fundamental theoretical inspiration for this undertaking. An introduction to the generic template for our dynamical models, namely the class of Galerkin spectral methods will naturally follow. Then we state the problem of buoyancy-driven convection in a fluid-saturated porous medium and review the phenomenological aspects thereof. This motivates the numerical investigation of the notion of the "minimal flow unit" in Chapter III. We show that the minimal flow unit may indeed serve as a dynamical unit to which the modeling can be reduced. In Chapter IV, we derive our new Galerkin method as well as the generic Fourier-Galerkin method, which we put to the test numerically in various dynamical regimes in Chapter V. Finally, Chapter VI presents our conclusions and closing remarks.

# CHAPTER II

# Background

Modeling, i.e., the process of constructing logical or mathematical frameworks that describe and predict natural phenomena in a sufficiently detailed and faithful fashion while being maximally parsimonious is arguably the cornerstone of the scientific method. In this sense, virtually any scientific theory, any equation describing a physical phenomenon, is ultimately a model: perhaps lacking in full accuracy and predictive power, but nevertheless capable of explaining and predicting to some finite extent.

The need for modeling may arise not only from our limited understanding of the fundamental physics, but also from the sheer mathematical complexity of physical systems even when the physics is exactly known. In classical physics, the Navier-Stokes equations describing the motion of fluids and other simplified variants are examples of such complex systems. Abstractly, they are represented by infinite-dimensional dynamical systems whose behavior can not be understood in complete analytical detail. Some such systems, however, possess a property that renders them in principle amenable to reduced modeling: their asymptotic dynamics are essentially finite-dimensional. It is to such systems that we now turn our attention.

## 2.1 Low-dimensional modeling

### 2.1.1 The global attractor

Driven dissipative nonlinear dynamical systems, including those evolving in ostensibly infinite dimensional phase spaces, often evolve on to invariant subsets after transients decay. These so-called global attractors contain the essential dynamical features of many complex systems. When applied to such systems, the theory of global attractors formalizes the idea that in certain cases, the asymptotic dynamics is essentially finite-dimensional, i.e., it displays enough regularity such that only a finite number of degrees of freedom are sufficient for a complete description of solutions contained in it.

The global attractor $\mathcal{A}$ can be defined as the maximal compact invariant set under the evolution semigroup defined on a Hilbert phase space. It exists if the semigroup is dissipative (as is the case in Navier-Stokes, Reaction-Diffusion and other equations) and provided there exists a compact absorbing set. It contains all complete bounded orbits and the unstable manifolds of all fixed points. The latter fact allows us to compute estimates the dimension of the attractor [2,3]. One can also prove that after long enough every orbit in the phase space will come to stay arbitrarily close to some trajectory on the global attractor, for an arbitrarily long time. This suggests that we must be able to approximate the asymptotic dynamics of the system with arbitrary precision by the dynamics restricted to the some superset of the global attractor. Temam [4] presents proofs of existence for the global attractors of several dissipative systems including reaction-diffusion equation, Navier-Stokes equations (in 2D), Rayleigh-Bénard (originally from [5]), and several dissipative wave equations. A proof of existence and bounds on the Hausdorff dimension of the global attractor of the porous-medium convection problem can be found in [6,7].

Although compact, in all but a few cases we are unable to prove that the global attractor is a smooth manifold. Therefore, dimensional bounds on the global attractor need to consider its fractal or Hausdorff dimensions. Furthermore, additional steps are required before we can formulate the asymptotic dynamics as the restriction of the semigroup on a finite-dimensional smooth manifold. Under certain conditions, the global attractor may be approximated by an "approximate manifold", a finite-dimensional sub-manifold of the phase space. The final step is the introduction of the stronger concept of the "inertial manifold" $\mathcal{M}$ defined as " A finite-dimensional Lipschitz manifold which is positively invariant and attracts all trajectories exponentially." [3, 8]

The existence of the inertial manifold is not trivial to prove or even known in many cases. For instance, while there exist proofs for many reaction-diffusion equations in one spatial dimension and in a rectangular domain in two dimensions, no proof is known for the generic two-dimensional geometry. Similarly, the existence of the inertial manifold for the Navier-Stokes equations in more than one dimension is an open problem (See Appendix A for details). In spite of this, the assumption of existence for typical dissipative systems is not too far-fetched in practice. In fact, in one way or another, it is implicit in any reduced modeling attempt. With this assumption, a natural reduced modeling strategy becomes readily available.

Consider an equation of the form

$$\frac{du}{dt} + Au = F(u) \quad u \in H \tag{2.1.1}$$

where $H$ is a Hilbert space, $A$ is a positive linear operator and $F$ is a Lipschitz function. Let $\{\psi_i,\ i = 1, 2, \cdots\}$ be the eigenfunctions of $A$ and let $P_N$ be the projection operator onto the

5

first $N$ eigenfunctions. Further, denote the orthogonal complement of $P_N H$ by $Q_N H$. We can expand $u$ in terms of the eigenfunctions: $u = \sum_1^\infty a_j \psi_j$ and $P_N u = \sum_1^N a_j \psi_j$. Clearly, $[A, P_N] = 0$. Then, by definition, for sufficiently large $N$, we expect that all "high" modal amplitudes, $a_j(t)$, $j > N$ be asymptotically "slaved" to the low modes. In other words, for $N$ large enough and $t \to \infty$, we expect that all solutions $u(t) = \sum_1^\infty a_j \psi_j$ be expressible as the sum of two terms: a linear combination of low modes, and a linear combination of high modes slaved to the low modes.

$$u(t) \to P_N\left(u(t)\right) + \phi\left(P_N(u(t))\right) \tag{2.1.2}$$

where $\phi : P_N H \to Q_N H$ is a Lipschitz function whose graph defines the inertial manifold. Then, in principle, the dynamics are determined fully by the low modes alone: the equation projected by $P_N$ together with the function $\phi$, fully characterize the asymptotic dynamics:

$$\frac{\mathrm{d}}{\mathrm{d}t} P_N u(t) + A P_N(u) = P_N F\left[P_N u + \phi(P_N(u))\right] \tag{2.1.3}$$

where we used the fact that $[A, P_N] = 0$. Thus, a set of $N$ ODEs will produce the asymptotic dynamics in full.

### 2.1.2 Galerkin projection

The functional form of the inertial manifold $\phi$ is almost never known explicitly. Consequently, we have to be content with the Galerkin truncation of the equation instead:

$$\frac{\mathrm{d}}{\mathrm{d}t} P_N u(t) + A P_N(u) = P_N F\left[P_N u\right]. \tag{2.1.4}$$

This is no longer the projection of the exact equation, but an equation for a low-mode approximation modified from the original equation to be entirely "blind" to the high modes. While the use of the eigenfunctions of the particular linear operator $A$ casts the Galerkin projected equation into an especially simple form, it is not technically necessary. One can perform the Galerkin projection onto any complete basis. In numerical analysis, this procedure constitutes the essence of Galerkin spectral methods sometimes used for efficient numerical solution of ODEs and PDEs, although less commonly than the other class, namely the pseudospectral collocation methods [9, 10]. In the theory of partial differential equations, it is particularly common to use the Fourier basis generically to produce finite truncations and generalize the obtained properties to the full PDE by proving convergence in the limit $N \rightarrow \infty$.

In practice, the choice of the basis is the defining feature of reduced-dimensional modeling strategies. Once projected onto the chosen basis, the equation assumes the form of a countable number of "modes" interacting with one another with various computable coupling coefficients. This presents a novel opportunity to gain insight into the physical processes involved: if chosen judiciously, the modes may be interpreted as representations of physical spatio-temporal "structures" that drive, inhibit or balance one another. A low-dimensional truncated model performing equally as well as a higher-dimensional model likely indicates a more efficient encoding of the dominant physical structures and their interaction processes into the selected basis functions. A physically motivated basis is likely to produce more efficient numerical methods, and conversely, an efficient method is likely a sign of a physically informative underlying basis.

### 2.1.3  Proper Orthogonal Decomposition

A traditional method for computing the "optimal" basis functions is the Proper Orthogonal Decomposition or POD. It is known in different contexts as the Karhunen-Loève decomposition, principal components analysis and singular value decomposition as well [11]. POD is an *a posteriori* statistical method requiring a sufficiently large ensemble of empirical data which must be obtained either experimentally or via direct numerical simulation. In essence, the method consists of statistically computing the principal axes of the "mass" of observed states residing in a Hilbert phase space. The ordered eigen-directions are then used as the new basis for the space.

More specifically, let $\{u_k\}_{k=1}^{\infty} \subset L^2(\Omega)$ be a representative ensemble of scalar functions, say, temperature fields in a thermal convection problem, obtained empirically. In order to find a complete set of normalized basis functions $\{\phi_j\}_{j=1}^{\infty}$ that maximally capture the "energy" or the $L^2$ norm of the $u_k$, we solve the following variational problem:

$$\max_{\phi \in L^2(\Omega)} \frac{\left\langle |(\phi, u)|^2 \right\rangle}{\|\phi\|^2} \tag{2.1.5}$$

where $\langle \cdots \rangle$ is the ensemble average over all $u_k$, $(\cdot, \cdot)$ denotes the inner product on $L^2(\Omega)$ and $\|\cdot\|$ is the $L^2$ norm. This leads to the maximization of

$$\left\langle |(\phi, u)|^2 \right\rangle - \lambda \left( \|\phi\|^2 - 1 \right). \tag{2.1.6}$$

The resulting Euler-Lagrange equation is the eigenvalue problem $\Re\phi = \lambda\phi$ [12] where

$$(\Re\phi)(x) = \int_\Omega \langle u(x)u^*(x') \rangle \, \phi(x') \, \mathrm{d}x' \tag{2.1.7}$$

8

with the kernel being simply the autocorrelation function of the empirical $u_k$. In [11], relatively lax necessary conditions are derived for $\mathfrak{R}$ to be a compact self-adjoint operator, in which case it possesses a complete eigenbasis and the expression in (2.1.5) has a maximum equal to the largest eigenvalue $\lambda_{\max}$.

Furthermore, given the spectral decomposition of a function in $\{u_k\}_{k=1}^{\infty}$ as $u(x) = \sum_1^{\infty} a_j \phi_j(x)$, one can show that

$$\langle a_j a_k^* \rangle = \delta_{jk} \lambda_j. \tag{2.1.8}$$

In other words, the modal amplitudes of the empirical functions in this basis are statistically uncorrelated. Additionally, $\lambda_j$ (which are necessarily positive) give the "average energy" in the $j$-th mode. For more details of the derivations see [11].

## 2.1.4 Discussion

The basis obtained as above can be proven to be optimal in the sense that the total ensemble-average energy in its first $N$ modes is greater that that of any other basis [11]. Thus, the use of POD constitutes a clear trade off: In order to attain energetic optimality, one has to measure a complete history of the dynamics and thus give up truly predictive power. The method serves as a powerful tool for empirical analysis of the dynamics but clearly fails to yield predictive a priori spectral methods.

Nevertheless, POD has been used to construct highly truncated models for a number of problems including turbulence coherent structures in Navier-Stokes [11] and turbulent Rayleigh-Bénard convection [13].

We, on the other hand, set out to derive *a priori* Galerkin methods for the problem of

porous-medium convection. The main goal of this dissertation is to explore how to compute *a priori* bases that are adapted specifically to the equations at hand and to the given parameter regime. The theoretical motivation also gives us hope of gaining insights into some of the coherent structures emerging in the porous-medium problem. We begin by stating the problem and reviewing the past phenomenological studies pertaining to porous-medium convection. The physical insights derived from the phenomenology, supplemented by our own direct numerical simulations, will guide us to assemble certain elements of our modeling strategy. The theoretical core of the method together with the complete derivations will follow that, and finally, numerical simulations of the reduced models constructed using the method will be studied extensively.

## 2.2   Porous-medium convection

The problem of buoyancy-driven convection in a fluid-saturated porous medium is a particularly simple yet elegant variant of Rayleigh-Bénard convection [14]. It has been successfully used to model geophysical phenomena where chemical or thermal inhomogeneity in fluids produces buoyancy and drives convection. Examples are geothermal reservoirs and carbon sequestration in saline aquifers [15]. Furthermore, despite its relative formal simplicity, it exhibits a range of increasingly complex behaviors as the main control parameter, the Rayleigh number ($Ra$) is increased. Although the flow eventually falls into a state of spatio-temporal chaos, it is nevertheless organized into recognizable large-scale structures that seem to follow a simple order, setting it apart from standard fully developed homogeneous isotropic turbulence as understood in the context of the Navier-Stokes equations [16]. Hence, the term turbulence, which we occasionally use to describe this regime, should be understood in a

qualified fashion.

It is this balance between mathematical tractability, dynamical complexity, and statistical spatio-temporal regularity which renders the problem of porous-medium convection especially appealing as a suitable model for the study of coherent structures and reduced modeling. There exists a sizable and growing literature on various experimental, computational and theoretical aspects of this problem on which parts of the present study are based [5,6,17–24].

## 2.2.1  Statement of the problem.

Consider the general three-dimensional geometry where a fluid saturated porous medium is placed in a rectangular box of height $h$, width $L$ and depth $d$. The state of the system is uniquely characterized by the temperature field $T(x,z,t)$ and the velocity field $\boldsymbol{u}(x,z,t) = u(x,z,t)\hat{\boldsymbol{i}} + v(x,z,t)\hat{\boldsymbol{j}} + w(x,z,t)\hat{\boldsymbol{k}}$.



Figure 2.2.1: The geometry of the three-dimensional setup and the boundary conditions.

We adopt the Darcy-Oberbeck-Boussinesq equations of motion [25]: the time evolution of the

temperature field is determined by an advection-diffusion equation which couples the temperature to a divergence-free velocity field. The momentum equation governing the evolution of the velocity field is given by Darcy's law in the Boussinesq approximation: the Navier-Stokes equation with internal forcing proportional to local velocity and an external forcing proportional to the local temperature to model thermal bouyancy. Darcy's law is a phenomenological equation describing the flow of a fluid in a porous medium. It states that the flux through a uniform block of a porous medium is proportional to the pressure drop across the medium. Equivalently, it implies that a fluid passing through a porous medium experiences a net resistance force per unit volume proportional to its velocity: $-\frac{\nu}{K}\boldsymbol{u}$. In modeling porous medium flow, this forcing replaces the viscous term $-\nu\Delta\boldsymbol{u}$ which in the Navier-Stokes equations represents fluid-fluid friction. Interestingly, the porous-medium forcing depends not only on $K$ (the square of the pore length scale) but also $\nu$, meaning it is still partially a product of fluid-fluid friction, albeit at the sub-pore length-scale. The Boussinesq approximation consists of neglecting density variations—and thus assuming incompressibility—except where they give rise to buoyancy: $g\alpha(T_{hot} - T_{cold})$. In the absence of any density gradient, buoyancy and therefore convection are impossible. Thus we have

$$T_t + \boldsymbol{u} \cdot \boldsymbol{\nabla}T = \kappa\Delta T \tag{2.2.1}$$

$$\boldsymbol{\nabla} \cdot \boldsymbol{u} = 0 \tag{2.2.2}$$

$$\boldsymbol{u}_t + \boldsymbol{u} \cdot \boldsymbol{\nabla}\boldsymbol{u} + \frac{\nu}{K}\boldsymbol{u} + \boldsymbol{\nabla}p = g\alpha(T_{hot} - T_{cold})\hat{\boldsymbol{k}}. \tag{2.2.3}$$

We consider a two-dimensional version of this problem where the fluid occupies a box of height $h$ and width $L$. This is not an entirely unphysical abstraction. In practice, the two-

12

dimensional version of the equations provide a good model for the experimentally realizable convection in the Hele-Shaw geometry where a an infinitesimally thin layer of fluid confined between two parallel vertical plates is studied.

The fluid is heated from below and cooled from above:

$$T(x, 0, t) = T_{hot} \quad \text{and} \quad T(x, h, t) = T_{cold} \quad \forall t > 0. \tag{2.2.4}$$

In the $x$ dimension, periodic boundary conditions are imposed for simplicity:

$$T(0, z, t) = T(L, z, t) \tag{2.2.5}$$

$$\boldsymbol{u}(0, z, t) = \boldsymbol{u}(L, z, t) \tag{2.2.6}$$

Further, we impose impenetrable boundary conditions at the top and bottom boundaries:

$$\boldsymbol{u} \cdot \hat{k} = w = 0 \quad \text{at} \quad z = 0, h. \tag{2.2.7}$$

Table 2.1 summarizes the physical parameters appearing in equations $(2.2.1) - (2.2.3)$.

| Parameter | Description |
|:---:|:---:|
| $\alpha$ | Thermal expansion coefficient |
| $K$ | Darcy Permeability coefficient |
| $h$ | Height of the layer |
| $\nu$ | Momentum diffusivity |
| $\kappa$ | Heat diffusivity |

Table 2.1: List of the physical parameters.

We choose the non-dimensionalization scheme used in [24] where time is measured in units of $h^2/\kappa$, length is measured in units of $h$, and temperature is measured in units of $(T_{hot} - T_{cold})$.

The dimensionless equations take a particularly simple form where all physical parameters are combined into two dimensionless parameters, namely the Rayleigh number and the Prandtl number $B^{-1}$ :

$$T_t + \boldsymbol{u} \cdot \boldsymbol{\nabla} T = \Delta T \tag{2.2.8}$$

$$\boldsymbol{\nabla} \cdot \boldsymbol{u} = 0 \tag{2.2.9}$$

$$B\left(\boldsymbol{u}_t + \boldsymbol{u} \cdot \boldsymbol{\nabla}\boldsymbol{u}\right) + \boldsymbol{u} + \boldsymbol{\nabla} p = Ra T \hat{\boldsymbol{k}} \tag{2.2.10}$$

where

$$Ra = \frac{g\alpha(T_{hot} - T_{cold})Kh}{\nu\kappa} \tag{2.2.11}$$

and

$$B^{-1} = \frac{\nu h^2}{\kappa K}. \tag{2.2.12}$$

The Rayleigh number is proportional to the overall temperature gradient $(T_{hot} - T_{cold})$ and is therefore a measure of the intensity of the thermal forcing. The Prandlt number is inversely proportional to the Darcy permeability $K$ which measures the square of the pore length scale. What we are concerned with is the infinite Prandtl number limit of the Darcy-Oberbeck-Boussinesq equations where the velocity $\boldsymbol{u}$ becomes slaved to the temperature $T$: it becomes a linear albeit non-local functional of $T$. Furthermore, the only nonlinearity is now in the coupling of the temperature with the velocity in the advection-diffusion equation. The equations in their final dimensionless form then reduce to

$$T_t + \boldsymbol{u} \cdot \boldsymbol{\nabla} T = \Delta T \tag{2.2.13}$$

$$\boldsymbol{\nabla} \cdot \boldsymbol{u} = 0 \quad \boldsymbol{u} = u\hat{\boldsymbol{i}} + w\hat{\boldsymbol{k}}. \tag{2.2.14}$$

14

$$\boldsymbol{u} + \boldsymbol{\nabla} p = RaT\hat{\boldsymbol{k}} \tag{2.2.15}$$

subject to boundary conditions

$$T(z = 0) = 1, \quad T(z = 1) = 0 \tag{2.2.16}$$

$$w(z = 0) = w(z = 1) = 0 \tag{2.2.17}$$

### 2.2.2   Nusselt number

The central emergent physical quantity of interest in our study of porous-medium convection is the dimensionless Nusselt number, denoted $Nu$ and defined as the spatially and long-time averaged vertical heat flux through the medium normalized by the heat flux due to conduction (diffusion) alone. In applications where the buoyancy originates from density gradients due to varying concentrations of different chemicals rather than temperature variations within the same material, for instance in the problem of CO2 sequestration [26], the Nusselt number measures mass transport rather than heat transport.

The various dynamical regimes encountered as $Ra$ is changed are conveniently characterized by the manner in which $Nu$ depends on $Ra$. For instance the onset of convection is (clearly) marked by the departure of $Nu$ from one, and the "turbulent" high-$Ra$ regime has come to be identified with the power law dependence of $Nu$ on $Ra$ in that regime.

In this study, we use the prediction of the Nusselt number as a measure of the "success" of our reduced models. In particular, by measuring the Nusselt number for models of increasing sizes (number of modes) at a given $Ra$ and tracking the trends thereof, we are able to gauge the rate and manner of convergence of our models to the exact system defined by the full

PDEs.

It is possible to derive a number of different expressions for $Nu$, all of which should yield identical results in the limit of long averaging time and high resolution. However, at severe truncations, as is the subject of our study, they need not be identical. Furthermore, due to issues of numerical stability, some may be more suitable for numerical evaluations than others. Here we derive all such expressions.

The heat flux density vector $\boldsymbol{J}$ (divided by density and the heat capacity) consists of a convective part and a diffusive part:

$$\boldsymbol{J} = \boldsymbol{u}T - \kappa \boldsymbol{\nabla} T. \tag{2.2.18}$$

The spatially averaged heat flux at a given time $t$ is then

$$Q(t) = \frac{1}{hL} \int_0^h \int_0^L \boldsymbol{J} \cdot \hat{\boldsymbol{k}} \, \mathrm{d}x \mathrm{d}z = \frac{1}{hL} \int_0^h \int_0^L wT \, \mathrm{d}x \mathrm{d}z - \frac{1}{hL} \int_0^h \int_0^L \kappa \partial_z T \, \mathrm{d}x \mathrm{d}z. \tag{2.2.19}$$

Denoting the area of the domain by $A = hL$ and defining the spatial and long- time average $\langle \cdot \rangle$ by

$$\langle \cdots \rangle = \lim_{t \to \infty} \frac{1}{t} \int_0^t \left[ \frac{1}{A} \int_0^h \int_0^L (\cdots) \, \mathrm{d}x \mathrm{d}z \right] \mathrm{d}\tau, \tag{2.2.20}$$

we obtain

$$\langle Q \rangle = \lim_{t \to \infty} \frac{1}{t} \int_0^t Q(t) \, \mathrm{d}t \tag{2.2.21}$$

$$= \langle wT \rangle + \frac{\kappa}{h^2} h \left[ T_{hot} - T_{cold} \right]. \tag{2.2.22}$$

16

The conductive heat flux can be easily shown to be equal to $\frac{\kappa}{h^2}h\left[T_{hot} - T_{cold}\right]$, equal to one In the dimensionless units. Thus the Nusselt number $\langle Q \rangle / \langle Q \rangle_{cond}$ is simply

$$Nu = 1 + \langle wT \rangle \tag{2.2.23}$$

in the dimensionless units. Henceforth, we will only work in the dimensionless units.

Another useful formulation is in terms of the Fourier transforms of the fields $w$ and $T$. Let $w_k(z)$ and $T_k(z)$ be the coefficients in the Fourier expansion of $w$ and $T$ in $x$ respectively:

$$T(x, z, t) = \sum_{k=-\infty}^{\infty} T_k(z, t)e^{\frac{2\pi i}{L}k} \tag{2.2.24}$$

$$w(x, z, t) = \sum_{k=-\infty}^{\infty} w_k(z, t)e^{\frac{2\pi i}{L}k}. \tag{2.2.25}$$

Then, by Parseval's theorem [27],

$$\frac{1}{L}\int_0^L wT \, \mathrm{d}x = \sum_{k=-\infty}^{\infty} \bar{w}_k(z, t)T_k(z, t) \tag{2.2.26}$$

so that

$$Nu = 1 + \lim_{t\to\infty} \frac{1}{t}\int_0^t \left[\sum_{k=-\infty}^{\infty}\int_0^1 \bar{w}_k(z, \tau)T_k(z, \tau) \, \mathrm{d}z\right] \mathrm{d}z. \tag{2.2.27}$$

This expression is particularly convenient in conjunction with the pseudospectral collocation methods of chapter III where the Fourier transforms of the fields are available and fast and efficient integration in the $z$ variable is possible on the Chebyshev grid.

17

## 2.2.3 Phenomenology

Porous-medium convection has been the subject of numerous experimental and numerical studies. Velocity field measurements are generally more difficult in the case of porous-medium convection than Rayleigh-Bénard in a pure fluid layer due to the presence of the porous medium. Thus, most experiments focus on measuring the total heat flux instead [17]. Using the minimally invasive technique of Magnetic Resonance Imaging (MRI) fairly detailed measurements of the velocity field have been made as well [18].

Moreover, the validity of Darcy's law requires a very small ratio of pore length scale to box size. Absent this condition, in particular at very high $Ra$ where the smallest wavelengths of the flow approach the characteristic pore size, the flow will depend on the microscopic details of the porous material and the model fails [18].

Complementary numerical experiments have revealed additional quantitative and qualitative details of the flow, providing us with a fairly complete knowledge of the various regimes from the onset of convection up to nearly $Ra \sim 10^5$, well into the turbulent regime [1, 19–21].

Energy stability analysis shows that at small enough $Ra$, the linear conduction solution $T(x,z) = 1 - z$ is absolutely stable, meaning any perturbation to this solution decays (Figure 2.2.2a). As $Ra$ is increased, the conduction solution becomes unstable, leading to steady convection in the form of a single pair of convective rolls (Figure 2.2.2b). Linear stability analysis allows us to compute the exact critical Rayleigh number $Ra_c$ as well as the horizontal wave length of the leading instability. For a box of aspect ratio 2, $Ra_c = 4\pi^2$. From onset up to $Ra \sim 380$, the steady convective rolls prevail until a Hopf bifurcation leads to a time-dependent flow where small-scale instabilities at the boundary layers are convected horizontally, forming periodic and quasi-periodic "traveling waves" [21]. At around $Ra \sim 700$,

(a) $Ra = 20$.

(b) $Ra = 300$

(c) $Ra = 700$

(d) $Ra = 2000$

Figure 2.2.2: Schematic depictions of the distinct dynamical regimes in porous-medium convection. The curves are isothermal lines in the plane. (a) The conduction regime. $T(x, z) = 1 - z$. (b) Steady convection in the form of one convective roll. (c) Periodic plume formation around the convective roll. (d) The single convective roll is lost and the high-$Ra$ regime has begun.

these instabilities begin to grow into recognizable "plumes" which are advected periodically around the large scale rolls (Figure 2.2.2c). Near $Ra = 1000$, chaotic solutions are found but the single pair of rolls stably remains the dominant structure until $Ra \sim 1250$ where it becomes unstable and the roll structure is lost (Figure 2.2.2d). This transition is marked by a sharp drop in the Nusselt number [19] and the replacement of the large convective rolls with chaotic "mega-plumes" [1, 19].

Perhaps the most important physical feature of the porous-medium flow at medium to high $Ra$ is the formation of boundary layers adjacent to the top and bottom boundaries. The formation of the boundary layers together with the rise of the plumes, mark the beginning of the so-called "scaling regime" where heuristic arguments predict a power-law dependence

of the Nusselt number upon $Ra$. Much theoretical, experimental and numerical effort has been dedicated to studying the scaling of the flow in this regime.

The so-called "marginal stability argument" was originally introduced by Malkus and developed by Howard for Rayleigh-Bénard convection [28–30]. The same argument may be used for porous-medium convection as well[1]. The outline of the argument is as follows: at high enough $Ra$ a distinct boundary layer is formed at the top and bottom boundaries. That is, in the horizontally averaged temperature profile virtually half of the total temperature drop $\frac{1}{2}(T_{hot} - T_{cold})$ occurs within a short distance of the boundaries while it barely varies in the bulk. Since the vertical velocity vanishes at the boundaries, the primary mode of heat transport in the boundary layer is diffusion. This links the total time averaged heat flux at the boundary and consequently the Nusselt number directly to the thickness of the boundary layer $\delta$. Finally, the thickness of the boundary layer may be linked to the Rayleigh number via a heuristic 'marginal stability' argument: the average thickness of the boundary layer is determined by a balance between the tendency of thinner layers to grow via diffusion on the one hand, and thicker layers to shrink down due to instabilities giving rise to convective plumes on the other. The boundary layer is thus "marginally stable", meaning that at each $Ra$, it assumes a thickness such that the effective $Ra$ of the boundary layer, i.e., the $Ra$ computed for a box of height $\delta$ and overall temperature difference $\frac{1}{2}(T_{hot} - T_{cold})$, is equal to the critical Rayleigh number $Ra_c$, which marks the onset of convection. Thus, $Nu$ and $Ra$ are linked via $\delta$ and the so-called Howard-Malkus-Kolmogorov-Spiegel (HMKS) scaling [29] for porous-medium convection is obtained:

$$\delta \sim Ra^{-1} \quad \text{and} \quad Nu \sim Ra. \tag{2.2.28}$$

---

[1]the difference being that $Ra \sim h$ instead of $Ra \sim h^3$ and a different scaling is obtained [31]

Although numerical studies have shown [19] that up to $Ra = 10^4$ the exponent is slightly less than one $\left(Nu \sim Ra^{0.9}\right)$, more recent numerical studies [1] show that at higher Rayleigh $\left(Ra \leq 40000\right)$, there is a significant trend toward linear scaling, suggesting that the Howard-Malkus-Kolmogorov-Spiegel scaling may in fact be exact asymptotically.

### 2.2.4    The minimal flow unit

Concerning low-dimensional modeling, perhaps the most important feature of the flow at high $Ra$ is the emergence of the "mega plumes"; persistent pairs of rising hot plumes and adjacent falling cold plumes stretching across the entire height of the domain, save for the increasingly thin boundary layers and relatively thicker but still considerably thin regions neighboring the boundary layers where instabilities grow into small "proto-plumes". In this regime, both our numerical results (see Chapter III) and those of Hewitt *et al* [1] suggest that the Nusselt number as well as the lateral length scale of the mega-plumes are independent of the aspect ratio $L$ at high enough $L$.

In other words, the mechanisms governing the local dynamics of the mega plumes and thus determining their size are blind to the overall extent of the domain. In fact, once the columnar structure settles into its asymptotic statistically stationary state, it appears that the pairs of plumes evolve more or less independently of one another (Figure 2.2.3). This suggests strongly that at very high $Ra$, a spatially "extensive" system evolves where certain measures of complexity including the number of dynamical degrees of freedom and the Kolmogorov-Sinai entropy grow linearly with the system size (in this case the aspect ratio $L$) [32]. Put differently, it is reasonable to assume that the approximate manifold of the full dynamics is roughly the direct sum of an integer number of smaller equivalent manifolds each describing

Figure 2.2.3: A snapshot of the temperature field at $Ra = 4 \times 10^4$ where the rising and falling "mega-plumes" span most of the height of the box, leaving only a thin proto-plume forming region and a much thinner boundary layer (too thin to be seen here) near the top and bottom.

the asymptotic dynamics of one of the statistically identical spatially localized subsets of the system, or the "minimal flow units".

This assumption paves the way for a significant reduction in size of dynamical models aimed at predicting intensive properties of the flow such as the vertical heat flux. Rather than wastefully resolving an entire box of aspect ratio 2 for instance, one may conceivably model only a single pair of plumes in an aspect ratio just large enough for the purpose.

To this end, one requires a theoretical or experimentally motivated prediction of the dependence of the size of the minimal flow unit on $Ra$. Hewitt $et$ $al$ [1] directly measured the size of the minimal flow unit (as defined here) by measuring the horizontal Fourier transform of the field (obtained from high-resolution direct numerical simulations with $L = 2$) at $z = \frac{1}{2}$ and

averaging the wave number with the largest amplitude. They observe that this wavenumber $k_{max}$ approximately satisfies a scaling law and measure the pre-factor:

$$k_{max} \simeq 0.5 Ra^{0.4} \tag{2.2.29}$$

or equivalently, $L \simeq 4\pi Ra^{-0.4}$. As of the date of this dissertation's writing, no theoretical argument predicting this 2/5 scaling has been presented.

A different and perhaps more fitting definition of the minimal unit for our purpose addresses the manner in which the dynamics changes as the aspect ratio varies from large to very small values. Such an approach has the potential to verify the physical assumption that more than being merely the visual quantum of flow, the minimal flow unit is indeed the smallest unit containing the essential dynamics. This will be our license to construct models for the minimal unit only, and assert that the results reflect those of the entire system in the appropriate sense. The following chapter describes a numerical investigation of the minimal flow unit in this pragmatic sense.

# CHAPTER III

# Direct numerical simulations

In this chapter, we explore the notion of the minimal flow unit numerically. As outlined in the previous chapter, the results will be applied to the construction of the reduced-order models developed and tested in Chapter IV and V. Our investigation consists of a series of direct numerical simulations (DNS) of porous-medium convection at a wide range of Rayleigh numbers and several aspect ratios where we measured the Nusselt number. The goal is to construct a complete picture of the exact aspect ratio-dependence of $Nu$ at a relatively large range of $Ra$ values in the turbulent regime. The trends revealed will lead us to formulate an operational definition of the minimal flow unit which we then use for reduced-dimensional modeling.

In what follows, we will describe the numerical methods employed in the direct numerical simulations and present the results.

## 3.1 Numerical methods

Our direct numerical simulations are based on a pseudospectral collocation method [9]. We write the equations of motion in terms of fluctuations about the conduction solution $1 - z$:

$$\partial_t \theta = -\boldsymbol{u}.\boldsymbol{\nabla}\theta + \Delta\theta - w \tag{3.1.1}$$

$$\Delta w = Ra\frac{\partial^2 \theta}{\partial x^2} \tag{3.1.2}$$

$$\boldsymbol{\nabla} \cdot \boldsymbol{u} = 0. \tag{3.1.3}$$

Therefore, $\theta(x, z, t)$ satisfies Dirichlet boundary conditions in $z$ and periodic boundary conditions in $x$. The box is disctretized using a regular equispaced grid in the $x$ dimension and the Chebyshev-Gauss-Lobatto grid [9] in the $z$ dimension.

$$x_j = j\frac{L}{N_x}, \ \ j = 0, 1, \cdots, N_x \tag{3.1.4}$$

$$z_j = \cos\left(j\frac{\pi}{N_z}\right), \ \ j = 0, 1, \cdots N_z. \tag{3.1.5}$$

Derivatives in $x$ are computed in the Fourier space using the Discrete Fourier Transform (DFT) [33] via Fast Fourier Transform (FFT). In the $z$ dimension, derivatives are computed using the Chebyshev differentiation matrix [33]. A numerical solution of the time-dependent equation involves the following at each time step:

1. Given the temperature field $\theta(x, z)$, solve the momentum equation (3.1.2) for $w$ in spectral space.

2. Solve the incompressibility equation (3.1.3) for $u$ in spectral space.

3. Form the right hand side of the advection-diffusion equation (3.1.1) in real space.

4. Perform time-stepping for $\theta$ using suitable methods for each of the terms.

The right hand side of (3.1.1) is divided into the linear (diffusion and $w$ ) and the nonlinear (advection) terms. All terms are computed in the Fourier space and then transformed back to the real space. This involves solving equations (3.1.3) and (3.1.2), the latter being the most computationally expensive step of all. An efficient method is employed for this task and will be elaborated below.

The time-stepping is performed in two separate half-steps. At each step, first the advection term is integrated using the 2nd order Adams-Bashforth method [34], which for the generic equation $\frac{\partial \theta}{\partial t} = F\{\theta\}$ and time step $h$ may be written as

$$\theta_{n+1} - \theta_n = \frac{h}{2}\left(3F\{\theta_n\} - F\{\theta_{n-1}\}\right). \tag{3.1.6}$$

The largest eigenvalue of the second order Chebyshev differentiation matrix grows quartically with the vertical resolution [33], rendering the time integration problem for the diffusion term ill-conditioned. Thus, an implicit method is required for the time-stepping of the diffusion term in order to ensure stability and we use the Crank-Nicolson method [34], also second order:

$$\theta_{n+1} - \theta_n = h\left[\frac{F\{\theta_{n+1}\} + F\{\theta_n\}}{2}\right]. \tag{3.1.7}$$

De-aliasing [9] is performed at each time step. A constant time-step $dt$ is used throughout each simulation due to the fact that the Adams-Bashforth method is a multi-step method. The multi-step Adams-Bashforth method allows us to achieve the desired order of accuracy with only one function evaluation per time step and is thus computationally cost effective.

The method is initialized by integrating over multiple fractional time-steps using the Euler method.

### 3.1.1 Solution of the momentum equation[2]

A pseudospectral solution of (3.1.2) requires the transformation of the equation into appropriate spectral spaces in both spatial dimensions. Taking the discrete Fourier transform of the equation in the $x$ direction (using for instance, FFT), we have

$$\left[D^2 - (nk)^2\right]\hat{w}_n(z) = -(nk)^2 Ra\hat{\theta}_n(z) \tag{3.1.8}$$

where $D = \partial_z$ and

$$\theta(x,z) = \sum_{n=-N_x/2+1}^{N_x/2} \hat{\theta}_n(z)e^{iknx} \quad, \quad k = \frac{2\pi}{L}. \tag{3.1.9}$$

Then, for each wavenumber $n$, we have a second order linear differential equation in the $z$ variable only, with a known right-hand-side. Since $z$ is discretized using the Chebyshev–Gauss–Lobatto grid, the most straightforward solution would be to use the real-space Chebyshev differentiation matrix and solve the resulting linear equation. However, the Chebyshev differentiation matrix and therefore the full LHS matrix are dense, and this method is not practical. We must therefore use the Chebyshev Integration Method [35]. The rest of this section describes the method for completeness, and closely follows the notation and exposition of [35].

The Chebyshev integration method is the method of choice for the solution of inhomogeneous

---

[2]Chebyshev integration method.

linear differential equations on the Lobatto grid. In general, suppose we would like to solve

$$L\psi = \sum_{k=0}^{n}(m_{n-k}(x)D^k)\psi = f(x) \ , \ \ x \in \Omega = (a,b) \tag{3.1.10}$$

subject to the constraints

$$\mathcal{T}\psi = c \tag{3.1.11}$$

where $m_k$ are rational functions of $x$, $D^k$ denotes $k$th order differentiation with respect to $x$ and $\mathcal{T}$ is a linear functional of rank $n$. Our equation is an example with $L = [D^2 - (nk)^2 I]$ and $f(z) = -(nk)^2 Ra\hat{\theta}_n(z)$.

In short, the Chebyshev Integration Method consists of solving for the highest order derivative of $\psi(z)$ rather than $\psi(z)$ itself, and doing so in the Chebyshev spectral space rather than the real space. Thanks to a particular property of the Chebyshev polynomials (and a number of other families of orthogonal polynomials), the resulting linear system will consist of a linear combination of various powers of the tridiagonal Chebyshev-basis-representation of the integration operator rather than the dense real-space Chebyshev differentiation matrix. The equation may then be solved using $O(N_z)$ operations rather than $O(N_z^2)$. In what follows, we discuss the technical details of the method.

Let $\{\phi_n(z), \ n = 0, 1, \cdots\}$ be a complete family of orthogonal polynomials that solve the singular Sturm-Liouville problem

$$\frac{\mathrm{d}}{\mathrm{d}x}\left(p(z)\frac{\mathrm{d}}{\mathrm{d}x}\phi_n(z)\right) + w(z)\phi_n(z) = 0 \quad x \in (a,b) \tag{3.1.12}$$

where the weight $w(x)$ is non-negative and $p(x) \to 0$ as $x \to a, b$.

Examples of such families include the Jacobi polynomials (and as a special case of that, the

28

Chebyshev polynomials), the Laguerre and the Hermite polynomials [35].

Let $Q_m^n = \text{span}\{\phi_j,\ m \le j \le n\}$. The Galerkin truncation of a function $f(z)$ in $Q_m^n$ is then

$$f(z) = \sum_{k=m}^{n} a_k \phi_k(z). \tag{3.1.13}$$

If $\phi_n$ are the Chebyshev polynomials, then the real-space representation of $f(z)$ on the $N$-point Lobatto grid can be used to efficiently compute its spectral representation in $Q_0^N$ , i.e., the coefficients $a_k$, $0 \le k \le N$, using Fast Fourier Transform. It is in this space that the equation may be solved efficiently. Let $\overline{f}$ be the vector of the Chebyshev spectral coefficients of the function $f(z)$, and let $\overline{L}$ and $\overline{D}$ be the Galerkin truncation of the linear operator $L$ and $D$ in the Chebyshev spectral space:

Then

$$\frac{1}{2}\overline{D}_{i,j} = \begin{cases} 0 & i \ge j \\ 0 & i < j,\ i+j \quad \text{even} \\ j & 0 < i < j,\ i+j \quad \text{odd} \\ \frac{j}{2} & i = 0,\ j \quad \text{odd}. \end{cases} \tag{3.1.14}$$

The recursion relation satisfied by the solutions of (3.1.12) [36] gives rise to "Integration Operators" that are banded; a property that is central to the efficiency of the method [35].

In particular, the integration operator in the Chebyshev basis is given by

$$
\overline{B} = \frac{1}{2} \begin{pmatrix}
0 & 0 & 0 & \cdots & 0 & \cdots & 0 \\
2 & 0 & -1 & \cdots & \cdots & \cdots & 0 \\
0 & 1/2 & 0 & -1/2 & \cdots & \cdots & 0 \\
0 & 0 & \ddots & \ddots & \ddots & \cdots & 0 \\
0 & 0 & 0 & 1/k & 0 & -1/k & \cdots
\end{pmatrix}.
\tag{3.1.15}
$$

$\overline{B}$ may be seen as the "inverse" of $\overline{D}$ by which we mean, $\overline{DB} = I_{Q_0^\infty} = I$. However, $BD = I_{Q_0^\infty} \neq I$. Also, in general, $\overline{D}^k \overline{B}^k = I$, but $\overline{B}^k \overline{D}^k \neq I$. The reason for that is the fact that differentiating a function $k$ times, erases all information about the first $k$ basis functions, and therefore, "integrating" the $k$-th derivative does not yield a unique result. However, if the domain of $\overline{D}$ and the range of $\overline{B}$ are properly restricted, then $\overline{B}$ truly becomes the inverse of $\overline{D}$. We will do this for the finite truncations of these operators. Note that $\overline{D} : Q_1^N \to Q_0^{N-1}$ and $\overline{D}^k : Q_k^N \to Q_0^{N-k}$. In order to construct the inverse of $\overline{D}^k$ in $Q_k^N$, we need two modifications to $\overline{B}^k$: we need to set the first $k$ rows of $\overline{B}^k$ to zero (thus setting its range to $Q_k^N$) and set its last $k$ columns to zero (thus changing its domain to $Q_0^{N-k}$). Let us call this modified version the "reduced" integration matrix and denote it by $\overline{B}_{[k]}^k$. Then, we have

$$
\overline{B}_{[k]}^k \overline{D}^k = I_{Q_k^N}.
\tag{3.1.16}
$$

Let us consider the porous-medium convection example:

$$
\left[ D^2 - \alpha \right] \psi(z) = f(z)
\tag{3.1.17}
$$

thus

$$\left[\overline{D}^2 - \alpha\right]\overline{\psi} = \overline{f} \tag{3.1.18}$$

The null-space of $\overline{L} = \left[\overline{D}^2 - \alpha\right]$, $\mathcal{N}(\overline{L})$ has dimension 2, reflecting the fact that the solution of $\overline{L}\overline{u} = 0$ involves two free parameters.

Write

$$\overline{\psi} = w + \psi_p \tag{3.1.19}$$

where $w \in \mathcal{N}(\overline{L})$ and $\psi_p \in Q_2^N$ is a particular solution: $\overline{L}\psi_p = \overline{f}$.

First, we solve for the particular solution by writing $\zeta = D^2\psi_p \in Q_0^{N-2}$ so that $\psi_p = \overline{B}_{[2]}^2\zeta \in Q_2^N$ is uniquely defined. Then the equation becomes

$$\left[I - \alpha\overline{B}_{[2]}^2\right]\zeta = \overline{f}. \tag{3.1.20}$$

Next, we find a basis for $\mathcal{N}(\overline{L})$ by writing

$$e_0 = \phi_0 + w_0 \tag{3.1.21}$$

$$e_1 = \phi_1 + w_1 \tag{3.1.22}$$

where $w_k \in Q_2^N$. Then, $\overline{L}e_k = 0$ yields

$$\overline{L}w_0 = -\overline{L}\phi_0 \tag{3.1.23}$$

$$\overline{L}w_1 = -\overline{L}\phi_1 \tag{3.1.24}$$

These equations need to be solved only once, and the result can be used for arbitrary right

hand sides ($\overline{f}$). Using this basis, we may now impose the boundary conditions

$$w = \sum_{k=0}^{1} a_k e_k \quad \Longrightarrow \quad \mathcal{T}\overline{\psi} = \mathcal{T}\psi_p + \sum_{k=0}^{1} a_k \mathcal{T} e_k = c. \tag{3.1.25}$$

Since $\mathcal{T}$ and the $e_k$ are known, we can compute $\mathcal{T}_{kl} = (\mathcal{T}e_k)_l$ , $l = 0, 1$. Therefore, to solve for the unknowns $a_k$ , $k = 0, 1$, we need to solve the $2 \times 2$ linear equation

$$\sum_{k=0}^{1} \mathcal{T}_{kl} a_k = c_l - (\mathcal{T}\psi_p)_l. \tag{3.1.26}$$

Having found both the particular solution and the general solution (by which the boundary conditions were imposed), we can now construct the full solution:

$$\overline{\psi} = w + \overline{B}_{[2]}^2 \zeta. \tag{3.1.27}$$

This last step involves a matrix multiplication, but $\overline{B}_{[2]}^2$ is also banded and thus the multiplication is an $O(N)$ operation.

## 3.1.2 Boundary conditions

In our problem the conditions enforced by $\mathcal{T}\overline{\psi} = c$ are homogeneous Dirichlet boundary conditions at the two limits of the domain. Writing the Chebyshev spectral expansion $\overline{\psi}(z) = \sum_{k=0}^{N} a_k T_n(z)$, $z \in [-1, 1]$ where $T_n$ is the $n$-th Chebyshev polynomial of the first kind [37] and $T_n(\cos(z)) = \cos(nz)$, we find that $\overline{\psi}(1) = a_0 + a_1 + a_2 + \cdots + a_N = 0$ and $\overline{\psi}(-1) = a_0 - a_1 + a_2 - \cdots \pm a_N = 0$. Equivalently, these conditions imply $a_0 + a_2 + a_4 + \cdots = 0$

and $a_1 + a_3 + a_5 + \cdots = 0$. We conclude that

$$
\begin{bmatrix}
1 & 0 & 1 & 0 & 1 & 0 & \cdots \\
0 & 1 & 0 & 1 & 0 & 1 & \cdots
\end{bmatrix}
\times \overline{\psi} =
\begin{bmatrix}
0 \\
0
\end{bmatrix},
\tag{3.1.28}
$$

or

$$
\mathcal{T} =
\begin{bmatrix}
1 & 0 & 1 & 0 & 1 & 0 & \cdots \\
0 & 1 & 0 & 1 & 0 & 1 & \cdots
\end{bmatrix}
\text{ and } c =
\begin{bmatrix}
0 \\
0
\end{bmatrix}.
\tag{3.1.29}
$$

## 3.2 Simulations

For the purposes of this chapter, we performed several independent simulations sampling the $Ra - L$ parameter space from $Ra \sim 10^3$ to $4 \times 10^4$ and $L \sim 0.05$ to 2. We measured $Nu^* = \langle Nu \rangle_L / \langle Nu \rangle_\infty$ where

$$
\langle Nu \rangle_L = \frac{1}{T} \int_0^T Nu(t) \, \mathrm{d}t \quad \text{for a box of width } L.
\tag{3.2.1}
$$

where $T$ is taken to be large enough so that the relative error in the measurement is no more than a few percent. For the range of $Ra$ studied, no significant variation in $\langle Nu \rangle_L$ was observed near $L = 2$ and therefore we assume that for our purposes $\langle Nu \rangle_\infty = \langle Nu \rangle_2$. The density of sample points is relatively sparse in the high $L$ region where there is little variation in $\langle Nu \rangle_L$ and more dense in the areas where high gradients are expected. For the most part, the sample points are logarithmically equi-spaced in $L$. A few points deemed to be outliers have been removed from the sample. We skipped points lying at high $Ra$ and high $L$ where $\langle Nu \rangle_L$ is expected to have converged to $\langle Nu \rangle_\infty$ and used the value 1 as the outcome. Figure 3.2.3(a) shows the sample space used for the study. Figure 3.2.3(b) is a

33

(a) $Ra = 700$     (b) $Ra = 11486$

Figure 3.2.1: Plots of $Ra$ vs $t$ at $Ra = 700$ and $Ra = 11486$. At high $Ra$, significantly higher temporal resolution is required for stability and accuracy.

color-coded map of $Nu^*$ in the $Ra-L$ plane with brighter colors representing higher values. A cubic interpolation of the discrete data points is used to smooth the final presentation. The contour lines represent the level sets for $Nu^* \in \{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 0.95\}$ from right to left. The level set for $Nu^* = 1$ has been omitted due to high levels of noise naturally occurring where $Nu^*$ plateaus.

Figure (3.2.4) shows snapshots of the temperature field obtained from the direct numerical simulations at Rayleigh numbers ranging from $10^3$ to $4 \times 10^4$. Using our large-aspect-ratio results—$L = 2$ at low $Ra$ and $L = 1$ at high $Ra$— juxtaposed with the results of Otero $et$ $al.$ [19], we are able to construct a full picture of $Nu$ vs $Ra$ for up to $Ra = 4 \times 10^4$. The results are plotted in figure (3.2.2).

Figure 3.2.2: Nusselt vs $Ra$ from direct numerical simulations. At high $Ra$, simulations are performed at aspect ratio 1 instead of 2.

## 3.3 Discussion

As evident from figure (3.2.4), the spacing of the mega-plumes indeed decreases with increasing $Ra$, qualitatively confirming the findings of Hewitt *et al* [1]. Our main goal however, is to propose a new definition for the minimal flow unit by studying the variations of the Nusselt number across the $Ra$-$L$ parameter space and compare it with that used by Hewitt *et al* [1].

We recognize the following regimes in the Nusselt landscape presented in (3.2.3)(b):

1. The conduction regime: regardless of $Ra$, at small enough $L$, convection becomes unsustainable and the flow settles to the conduction solution and $Nu = 1$. (Note that $Nu^*$ is small although non-zero). The onset of convection occurs at different

(a)



(b)

Figure 3.2.3: (a) A representation of the points sampled in the $Ra$-$L$ plane to produce the next plot. The blank area at high $Ra$ and $L$ was not sampled since physical arguments and numerical indications from the rest of the sample suggest that $Nu$ is essentially equal to $Nu^*$ in that region.

(b) A color-coded map of $Nu^* = \langle Nu \rangle_L / \langle Nu \rangle_\infty$ where $\langle Nu \rangle_\infty$ is the infinite-aspect-ratio Nusselt. The contour lines represent the level sets for $0.1, 0.2, \cdots, 0.9, 0.95$ from right to left. At small aspect ratios, eventually convection becomes unsustainable and $Nu \to 0$. The solid line scale corresponds to the minimal flow unit measurements of Hewitt *et al* [1].

36

(a) $Ra = 11486$    (b) $Ra = 13195$    (c) $Ra = 15157$

(d) $Ra = 17411$    (e) $Ra = 20000$    (f) $Ra = 22973$

(g) $Ra = 26390$    (h) $Ra = 30314$    (i) $Ra = 34822$

(j) $Ra = 40000$    (k) $Ra = 45947$

Figure 3.2.4: Snapshots of the temperature field obtained from direct numerical simulations at $L = 1$.

aspect ratios depending on $Ra$ and the precise critical $L$ may be computed using linear stability theory as follows. As derived (for example) in Doering & Constantin [24] the eigenvalues of the linear stability problem for a box of height 1 are given by

$$\lambda_{m,k} = k^2 + m^2\pi^2 - \frac{Rak^2}{k^2 + m^2\pi^2}. \tag{3.3.1}$$

where $k$ is the horizontal wave number and $m\pi$ is the vertical wave number ($m = 1, 2, 3, \cdots$). To investigate the onset of instability, we consider the lowest of the branches: $m = 1$. It is easily found that above the critical Rayleigh number $Ra^* = 4\pi^2$, $\lambda_{1,k} > 0$ for some open interval of $k$. For a periodic box of finite size however, instability requires that a quantized horizontal wave number $n = kL/(2\pi)$ fall within that interval. Therefore, the smallest $L$ such that an unstable mode exists corresponds to that for which the wave number with $n = 1$ coincides with the larger of the two roots of $\lambda_{1,k}$. Thus,
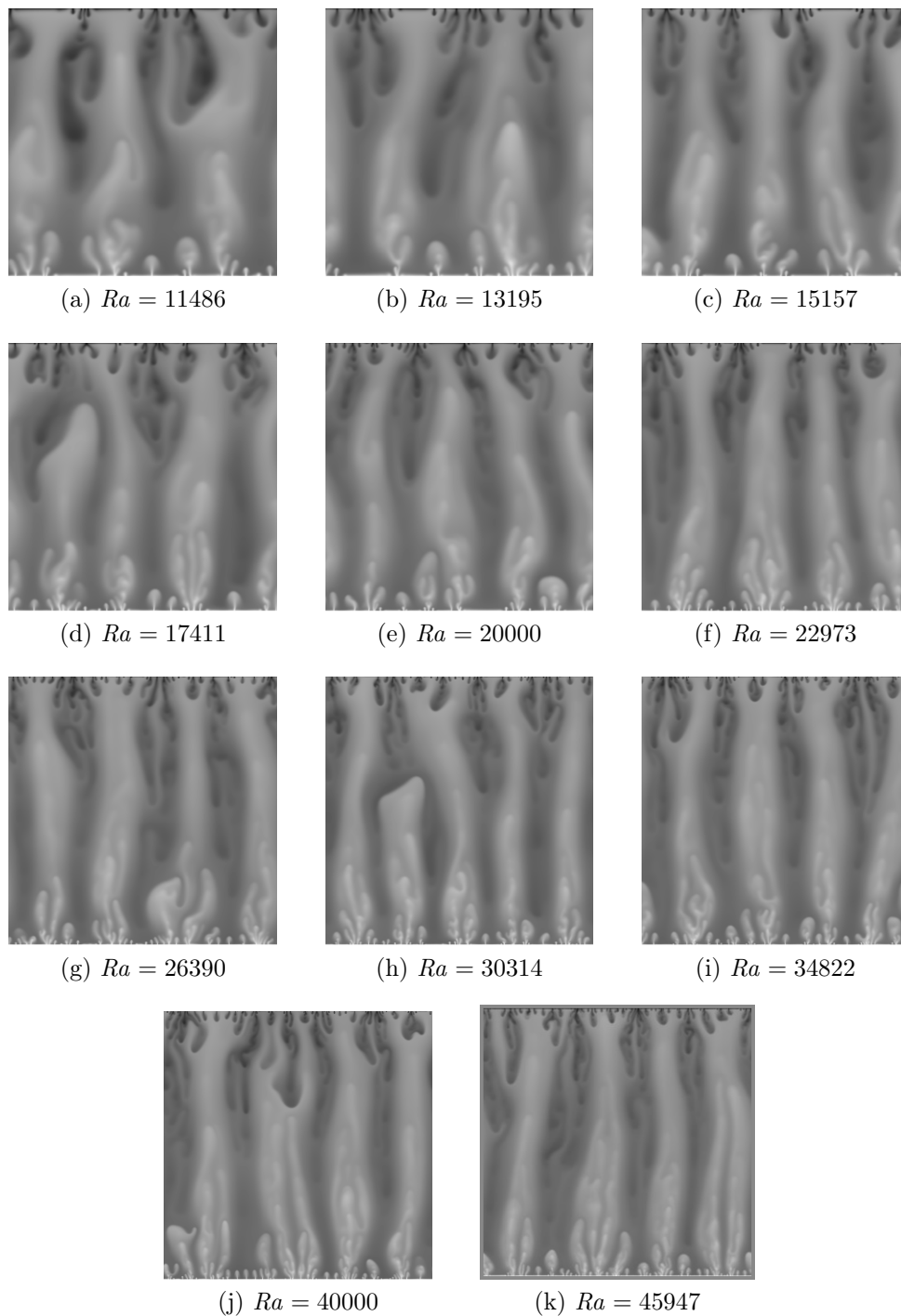
$$L_{\min} = \frac{4\pi}{Ra^{1/2} + \sqrt{Ra - 4\pi^2}} \sim 2\pi Ra^{-1/2} \quad \text{as} \quad Ra \to \infty. \tag{3.3.2}$$

This line indeed delineates the boundary of the dark region at the bottom right of figure 3.2.3 (b).

2. The large-aspect ratio regime. At any $Ra$, at large enough $L$, eventually $Nu$ becomes independent of $L$ and $Nu^*$ plateaus to 1. This is to be expected physically since in the turbulent regime, the flow is organized in apparently independent and statistically similar "cells", each consisting of a rising column and an adjacent falling column. Our result clearly confirms this at low to medium $Ra$ values. At high $Ra$ ($10^4 - 4 \times 10^4$), we have not performed high-$L$ simulations. However, the trend toward the $Nu^* = 1$ plateau is clearly visible from our low-$L$ simulations, and is consistent with results from

lower $Ra$.

3. The intermediate regime. Between the onset of convection and the high-$L$ plateau, convection is sustained but not fully actualized due to small box size, and $Nu^* < 1$. We interpret the boundary between this region and the $L$-independent plateau as the minimal box size that sustains a fully actualized convection cell.

At each $Ra$, the solid line in 3.2.3 (b) follows the scaling suggested by the time-averaged width of a mega-plume pair or a single convective cell measured by Hewitt $et$ $al$ [1].

$$k = 0.5 \, Ra^{\frac{2}{5}} \tag{3.3.3}$$

We believe that this line is remarkably close to the boundary between the $Nu^* = 1$ plateau and the intermediate regime. For reference, the blue line represents $k \sim Ra^{-1/2}$. Although not fully resolved at high $Ra$, the level sets of $Nu^*$ in 3.2.3 (b) seem to follow the $Ra^{\frac{2}{5}}$ scaling more closely.

Note that any coincidence between the time-averaged convection cell width in a large-aspect ratio box and our measurement of the minimal aspect ratio sustaining $Nu^* = 1$ is quite non-trivial. The former corresponds to the length scale naturally selected by the dynamics independently of any geometric constraints, while the latter demonstrates the behavior of the flow in response to imposed size constraints.

If true as we suspect, the coincidence of the two provides new insights into the dynamics: it suggests that even when the width of the box is infinite, the convective cells asymptotically self-organize into the smallest units that are just large enough to allow the maximum heat transport to be realized.

Henceforth, we adopt the scaling $k = 0.5\,Ra^{\frac{2}{5}}$ or equivalently $L = 4\pi Ra^{-\frac{2}{5}}$ as our operational definition of the minimal flow unit width. In some of the low-dimensional models we will derive and investigate in the coming chapters, we will use this aspect ratio for high-$Ra$ simulations.

# CHAPTER IV

# Galerkin methods from upper-bound theory

In this chapter, we derive Galerkin methods for the solution of porous-medium convection. First we derive a method based on the Fourier basis in both dimensions, henceforth referred to as the Fourier-Galerkin (FG) method. Using the Fourier basis as the most "neutral" basis possible provides us with a point of reference for future comparisons.

Then, we proceed to derive our adapted Galerkin method (Nonlinear-Galerkin method or NG), using a basis computed numerically at each Rayleigh number. This method has its roots in upper-bound theory [22–24, 38] and energy-stability theory which will be explained first. The results of extensive numerical studies of these methods will be presented in the next chapter.

## 4.1  The Fourier-Galerkin method

In order to derive the Fourier-Galerkin method, we project the evolution equation of the temperature field onto the two-dimensional Fourier basis satisfying the boundary conditions.

To simplify the boundary conditions, we rewrite the equations in terms of the fluctuations $\theta(x, z, t)$ about the conduction solution: $T(x, z, t) = (1 - z) + \theta(x, z, t)$. We will then have periodic boundary conditions in $x$ and Dirichlet boundary conditions in $z$. The temperature basis will consist of the functions $\phi_{mn}(x, z) = e^{inkx} \sin(mk'z)$, $m = 1, 2, \cdots \infty$ and $n = -\infty, \cdots +\infty$, $k = 2\pi/L$, $k' = \Gamma k$ where $\Gamma$ is the aspect ratio of the box. The no-flux boundary condition at the vertical boundaries and the incompressibility condition will dictate the basis to be used for the expansion of the velocity fields

$$\theta(x, z, t) = \sum_{n=-\infty}^{\infty} \sum_{m=1}^{\infty} a_{mn}(t) e^{inkx} \sin(mk'z) \tag{4.1.1}$$

$$u(x, z, t) = \sum_{n=-\infty}^{\infty} \sum_{m=1}^{\infty} b_{mn}(t) e^{inkx} \cos(mk'z) \tag{4.1.2}$$

$$w(x, z, t) = \sum_{n=-\infty}^{\infty} \sum_{m=1}^{\infty} c_{mn}(t) e^{inkx} \sin(mk'z). \tag{4.1.3}$$

Two constraints relate the modal amplitudes $a_{mn}(t)$, $b_{mn}(t)$, $c_{mn}(t)$ :

1. Incompressibility: $\partial_x u + \partial_z w = 0$. This implies

$$\sum_{n=-\infty}^{\infty} \sum_{m=1}^{\infty} (ink) b_{mn}(t) e^{inkx} \cos(mk'z) = \sum_{n=-\infty}^{\infty} \sum_{m=1}^{\infty} (mk') c_{mn}(t) e^{inkx} \sin(mk'z) \tag{4.1.4}$$

or

$$(ink) b_{mn}(t) = (mk') c_{mn}(t). \tag{4.1.5}$$

2. The momentum equation: $\Delta w + \text{Ra}(D^2 - \Delta)\theta = 0$. This similarly yields

$$\left[ (nk)^2 + (mk')^2 \right] c_{mn}(t) = Ra(nk)^2 a_{mn}(t). \tag{4.1.6}$$

42

Using these constraints we may now perform the full Galerkin projection of the time evolution equation

$$\partial_t \theta + u \partial_x \theta + w \partial_z \theta = w + \left( \partial_x^2 + \partial_z^2 \right) \theta. \tag{4.1.7}$$

The various terms can be expanded as follows:

$$\partial_t \theta = \sum_{n'=-\infty}^{\infty} \sum_{m'=1}^{\infty} \dot{a}_{m'n'}(t) e^{in'kx} \sin(mk'z) \tag{4.1.8}$$

$$u \partial_x \theta = \sum_{m',n'} \sum_{p,q} (in'k) a_{m'n'} b_{pq} e^{ikx(n'+q)} \sin(k'm'z) \cos(k'pz) \tag{4.1.9}$$

$$w \partial_z \theta = \sum_{m',n'} \sum_{p,q} (k'm') a_{m'n'} c_{pq} e^{ikx(n'+q)} \sin(k'pz) \cos(k'm'z) \tag{4.1.10}$$

$$\left( \partial_x^2 + \partial_z^2 \right) \theta = \sum_{n'=-\infty}^{\infty} \sum_{m'=1}^{\infty} \left[ -(kn')^2 - (k'm')^2 \right] a_{m'n'}(t) e^{in'kx} \sin(m'k'z). \tag{4.1.11}$$

Using these expansions, we may compute the projections of both sides of the equation onto any given basis function $\phi_{mn} = e^{inkx} \sin(mk'z)$ by taking the inner product

$$(\phi_{mn}, f) = \frac{2}{L} \int_0^1 \int_0^L e^{-inkx} \sin(mk'z) f(x,z) \mathrm{d}x \mathrm{d}z.$$

This yields

$$(\phi_{mn}, \partial_t \theta) = \dot{a}_{mn} \tag{4.1.12}$$

$$(\phi_{mn}, w) = \frac{Ra(nk)^2}{[(nk)^2 + (mk')^2]} a_{mn} \tag{4.1.13}$$

$$
\begin{aligned}
(\phi_{mn}, u\partial_x\theta) \;=\;& \frac{2}{L}\sum_{m',n'}\sum_{p,q}(in'k)a_{m'n'}b_{pq}e^{ikx(n'+q-n)} \\
& \times \int \sin(k'm'z)\cos(k'pz)\sin(mk'z)\mathrm{d}x\mathrm{d}z \qquad (4.1.14)\\
=\;& 2\sum_{m',p,q}(i(n-q)k)a_{m'(n-q)}b_{pq} \\
& \times \left[\frac{1}{4}\delta_{m',m-p}+\frac{1}{4}\delta_{m',m+p}-\frac{1}{4}\delta_{m',p-m}-\frac{1}{4}\delta_{m',-m-p}\right] \qquad (4.1.15)
\end{aligned}
$$

$$
\begin{aligned}
(\phi_{mn}, w\partial_z\theta) \;=\;& \frac{2}{L}\sum_{m',n'}\sum_{p,q}(m'k')a_{m'n'}c_{pq}e^{ikx(n'+q-n)} \\
& \times \int \sin(k'pz)\cos(k'm'z)\sin(mk'z)\mathrm{d}x\mathrm{d}z \qquad (4.1.16)\\
=\;& 2\sum_{m',p,q}(m'k')a_{m'(n-q)}c_{pq} \\
& \times \left[\frac{1}{4}\delta_{p,m-m'}+\frac{1}{4}\delta_{p,m+m'}-\frac{1}{4}\delta_{p,m'-m}-\frac{1}{4}\delta_{p,-m'-m}\right] \qquad (4.1.17)
\end{aligned}
$$

$$
\left(\phi_{mn}, \left[\partial_x^2+\partial_z^2\right]\theta\right) = -\left[(kn)^2+(k'm)^2\right]a_{mn} \qquad (4.1.18)
$$

where we have used the fact that

$$
\int_0^1 \sin(\pi n'z)\sin(\pi nz)\cos(\pi qz)\mathrm{d}z = 
\begin{cases}
\frac{1}{4} & n'=n-q \ \text{ or } \ n'=n+q \\[2mm]
-\frac{1}{4} & n'=q-n \ \text{ or } \ n'=-(n+q) \\[2mm]
0 & \text{otherwise.}
\end{cases}
\qquad (4.1.19)
$$

Reducing the sums over $m'$, we arrive at

$$
(\phi_{mn}, u\partial_x\theta) = \frac{1}{2}\sum_{p,q}i(n-q)kb_{pq}\left[a_{(m-p)(n-q)}+a_{(m+p)(n-q)}-a_{(p-m)(n-q)}-a_{(-m-p)(n-q)}\right]
$$
$$
(4.1.20)
$$

$$(\phi_{mn}, w\partial_z\theta) = \frac{1}{2}\sum_{p,q} c_{pq} \left[ k'(m-p)a_{(m-p)(n-q)} + k'(p-m)a_{(m+p)(n-q)} - k'(p-m)a_{(p-m)(n-q)} \right]$$

$$(4.1.21)$$

All terms put together, we obtain the following set of ordinary differential equations:

$$\dot{a}_{mn} = \mu_{mn}a_{mn} + \sum_{p,q}\alpha_{pq}^{mn}a_{pq}a_{(m-p)(n-q)} + \sum_{p,q}\beta_{pq}^{mn}a_{pq}a_{(m+p)(n-q)} + \sum_{p,q}\gamma_{pq}^{mn}a_{pq}a_{(p-m)(n-q)}$$

$$m = 1, 2, \cdots \quad \text{and} \quad n = -\infty, \cdots, \infty \qquad (4.1.22)$$

where

$$\alpha_{pq}^{mn} = \frac{1}{2}kk'\Gamma_{pq}\left[p(n-q) - q(m-p)\right] \qquad (4.1.23)$$

$$\beta_{pq}^{mn} = \frac{1}{2}kk'\Gamma_{pq}\left[p(n-q) - q(m+p)\right] \qquad (4.1.24)$$

$$\gamma_{pq}^{mn} = -\alpha_{pq}^{mn} \qquad (4.1.25)$$

$$\Gamma_{pq} = \frac{\text{Ra}(kq)}{(kq)^2 + (k'p)^2} \qquad (4.1.26)$$

and the linear coefficient is given by

$$\mu_{mn} = \frac{\text{Ra}(kn)^2}{(kn)^2 + (k'm)^2} - (kn)^2 - (k'm)^2. \qquad (4.1.27)$$

We have transformed the partial differential equations into a countably infinite set of ordinary differential equations. Given any finite subset of

$$A = \{(m_1, n_1), (m_2, n_2), \cdots, (m_N, n_N)\} \cup \{(m_1, -n_1), (m_2, -n_2), \cdots, (m_N, -n_N)\} \quad (4.1.28)$$

(symmetrized about zero in the second index to satisfy the reality condition), the corre-

sponding Galerkin approximation of the equations will be

$$
\begin{aligned}
\dot{a}_{mn} \;=\;\; & \mu_{mn} a_{mn} \\
& + \sum_{\substack{(p,q)\in A \\ (m-p,,n-q)\in A}} \alpha^{mn}_{pq} a_{pq} a_{(m-p)(n-q)} \\
& + \sum_{\substack{(p,q)\in A \\ (m+p,n-q)\in A}} \beta^{mn}_{pq} a_{pq} a_{(m+p)(n-q)} \\
& + \sum_{\substack{(p,q)\in A \\ (p-m,n-q)\in A}} \gamma^{mn}_{pq} a_{pq} a_{(p-m)(n-q)} \\
& (m,n) \in A
\end{aligned}
\tag{4.1.29}
$$

This finite set of ordinary differential equations may now be analyzed using the standard techniques of finite-dimensional dynamical systems.

## 4.2   The Nonlinear-Galerkin method

While standard, the use of the Fourier basis as demonstrated above fails to take advantage of the distinctive spatio-temporal regularities, or the coherent structures of the system in question. We may think of the Fourier basis as the most "neutral" of all bases, capable only of encoding length scale in a "blind" and spatially uniform manner. It fails for instance, to efficiently represent spatially localized structures such as the increasingly thin boundary layers encountered in Rayleigh-Bénard or porous-medium convection.

We ask if one can derive a basis from the particular equations of motion at hand, such that the basis itself reflects the dominant spatial features of the motion in some sense, and the Galerkin-projected dynamical systems capture the dynamics as modal interactions among a

46

subset of modes minimal in size and maximal in qualitative representation of spatial features of the motion.

Our strategy is inspired by the background field method [22–24, 38] where the temperature field is decomposed into a mean background profile and a fluctuation field about the background profile. Various properties of the fluctuation field may be controlled by the choice of the background profile and thus the analysis of the fluctuation field may be cast into the most suitable form– depending on one's goals– by solving appropriate variational problems involving the background profile.

In particular, we solve a variational problem yielding a background profile that captures the boundary layer structure of the flow and produces a good first approximation for the Nusselt number. At the same time, the solution to the variational problem produces a basis with remarkable properties, suggesting it could be used for efficient Galerkin methods.

### 4.2.1  Upper-bound theory

Let $\tau(z)$ be a smooth function satisfying the temperature boundary conditions $\tau(0) = 1$, $\tau(1) = 0$. Write $T(x, z, t) = \tau(z) + \theta(x, z, t)$ so that $\theta(x, z, t)$ is the time-dependent fluctuation about the "background profile" $\tau(z)$. Note that $\theta(x, z, t)$ now vanishes at $z = 0$ and $z = 1$, and satisfies periodic boundary conditions in $x$ similar to $T$. The equations of motion may now be written in terms of $\theta$ :

$$\partial_t \theta + \boldsymbol{u}.\boldsymbol{\nabla}\theta \;=\; \Delta\theta + \tau'' - w\tau' \tag{4.2.1}$$

$$\Delta w \;=\; Ra\frac{\partial^2\theta}{\partial x^2} \tag{4.2.2}$$

$$\boldsymbol{\nabla}\cdot\boldsymbol{u} \;=\; 0. \tag{4.2.3}$$

Using this background decomposition, Doering and Constantin [24] derived rigorous upper bounds for the Nusselt number as a function of $Ra$. They showed that provided the functional

$$\mathscr{H}_a\{\theta\} = \left\langle a\,|\nabla\theta|^2 + \theta w \tau'(z)\right\rangle \tag{4.2.4}$$

is positive semi-definite, we have the following rigorous upper bound for $Nu$ :

$$Nu - 1 \le \frac{1}{4a(1-a)}\left[\int_0^1 (\tau'(z))^2 dz - 1\right], \tag{4.2.5}$$

where $a \in (0,1)$ and $\langle\cdot\rangle = \int_0^1\int_0^L(\cdot)\mathrm{d}x\mathrm{d}z$. For $a = \frac{1}{2}$, the upper bound is precisely the Nusselt number produced by $\tau(z)$ if it were the full temperature field. Figure 4.2.3 shows the optimal upper bound alongside "exact" values of $Nu$ obtained from direct numerical simulations.

Noting that $w(x,z)$ is a linear functional of $\theta$, we see that $\mathscr{H}_a\{\theta\}$ is indeed a quadratic form in terms of $\theta$ :

$$\mathscr{H}_a\{\theta\} = \int_0^1\int_0^L \left[a\,|\nabla\theta|^2 + \theta w \tau'(z)\right]\mathrm{d}x\mathrm{d}z \tag{4.2.6}$$

$$= \int_0^1\int_0^L \left[-a\theta\Delta\theta + \theta\tau'w[\theta]\right]\mathrm{d}x\mathrm{d}z. \tag{4.2.7}$$

This may be rewritten as

$$\mathscr{H}_a\{\theta\} = (\theta, \mathscr{L}\{\theta\}) \tag{4.2.8}$$

where $(f,g) = \int_0^1\int_0^L f(x,z)g(x,z)\mathrm{d}x\mathrm{d}z$ indicates the inner product and $\mathscr{L} = -a\Delta + \tau'w$ is a linear operator acting on $\theta$. It is clear that the condition $\mathscr{H}_a\{\theta\} \ge 0$ is equivalent to a spectral constraint on $\mathscr{L}$. Let $\mathsf{S}$ and $\mathsf{A}$ be the symmetric and antisymmetric parts of $\mathscr{L}$ respectively:

$$\mathscr{L} = \mathsf{A} + \mathsf{S}. \tag{4.2.9}$$

Then, $\mathscr{H}_a\{\theta\} = (\theta, (\mathsf{A} + \mathsf{S})\{\theta\}) = (\theta, \mathsf{S}\{\theta\})$ and the spectral constraint may be equivalently applied to $\mathsf{S}$. Thus, we have identified a symmetric linear operator, depending parametrically on $\tau$, whose spectrum determines if $\tau$ produces an upper bound on $Nu$.

## 4.2.2 Energy-stability

We may interpret the spectral constraint in terms of energy-stability as follows: multiplying (4.2.1) by $\theta$ and integrating, we obtain

$$\frac{1}{2}\frac{\mathrm{d}}{\mathrm{d}t}\|\theta\|_2^2 = \int \theta(\Delta\theta - w\tau')\mathrm{d}x\mathrm{d}z + \int \theta\tau''\mathrm{d}x\mathrm{d}z \tag{4.2.10}$$

where $\|\cdot\|_2$ indicates the $L_2$ norm over the $x$-$z$ domain and the term $\int \theta\boldsymbol{u}\cdot\boldsymbol{\nabla}\theta\,\mathrm{d}x\mathrm{d}z$ vanishes due to incompressibility and the boundary conditions. It is evident that

$$\frac{1}{2}\frac{\mathrm{d}}{\mathrm{d}t}\|\theta\|_2^2 = -\mathscr{H}_1\{\theta\} + \int \theta\tau''\mathrm{d}x\mathrm{d}z. \tag{4.2.11}$$

Energy stability then requires that the sum of the two terms on the right hand side be negative definite. Were it not for the second term, this condition would have resembled the spectral constraint necessary for the upper-bound argument. It would not have produced an upper bound at $Ra$, since $a \in (0,1)$ in (4.2.5). However, it would have produced an upper bound at a lower $Ra$ such as $Ra/2$ since $\mathscr{H}_1\{\theta\}$ at $Ra$ is proportional to $\mathscr{H}_{\frac{1}{2}}\{\theta\}$ at $Ra/2$ due to the linear relationship between $w$ and $\theta$ (4.2.2). Consequently, we conclude that if $\tau$ is chosen such that the constraint $\mathscr{H}_1\{\theta\} \geq 0$ is satisfied, then $\tau$ produces a Nusselt upper

49

bound at $Ra/2$ (which may be optimized by varying $\tau$) while isolating the destabilizing factors in one "forcing term" ($\int \theta \tau'' \mathrm{d}x\mathrm{d}z$). The upper bound $\tau$ produced for $Ra/2$ in this manner proves to be a surprisingly good estimate for $Nu$ at $Ra$ as the upper bound plot illustrates.

## 4.2.3 Modeling strategy

Based on the arguments discussed above, we propose the following modeling strategy: We compute an "optimal" background profile $\tau(z)$ by numerically solving the variational problem

$$Nu = \inf_{\substack{\tau(0)=1 \\ \tau(1)=0}} \left\{ \int_0^1 \tau'^2(z)\, \mathrm{d}z \,|\, \mathscr{H}_1\{\theta\} \geq 0 \right\}. \tag{4.2.12}$$

This is a variational problem subject to a spectral constraint on the self-adjoint operator inside the quadratic form $\mathscr{H}_1\{\theta\}$. Thus, by solving this problem, we accomplish a number of goals:

1. We find an "optimal" background profile $\tau(z)$ which yields a good first approximation to the Nusselt number.

2. We find a complete orthogonal eigenbasis and the associated positive semi-definite spectrum which we can then use for the construction of Galerkin spectral methods. The linear modal coefficients in the resulting models turn out to be negatively proportional to the eigenvalues found, and are therefore all either zero or negative. A number of "marginally stable modes" will emerge from the solution which may indicate dominant dynamical length scales.

50

3. We successfully separate the forcing terms from the linear and non-linear modal inter-
   actions in the ODEs that constitute our Galerkin methods. This offers more versatility
   in choosing appropriate truncations while addressing numerical stability.

Next, we will formulate the eigenvalue problem and derive the Galerkin methods based on
the resulting eigenbasis.

## 4.2.4   The eigenvalue problem

For a fixed background profile $\tau$, the spectral condition is essentially a constrained variational
problem where both (4.2.2) and (4.2.3) have to be enforced implicitly. In order to derive
an explicit eigenvalue problem, we solve the variational problem, enforcing the constraints
using Lagrange multipliers.

We demand that the infimum of the $\mathscr{H}_1\{\theta\}$ over all smooth $\theta$ be non-negative. Since $\mathscr{H}_1\{\theta\}$
is a quadratic form in $\theta$, a necessary and sufficient condition for this is that the infimum be
taken over all smooth $\theta$ of some fixed $L_2$ norm(say, 1), or equivalently that we compute the
infimum of $\mathscr{H}_1\{\theta\}$ normalized by $\|\theta\|_2^2$. We proceed with the former condition for reasons
that will become apparent shortly.

$$\lambda_{min} = \inf_\theta \int \left( |\nabla\theta|^2 + \theta w \tau' \right) \, \mathrm{d}x\mathrm{d}z, \quad \int \theta^2 \mathrm{d}x\mathrm{d}z = 1. \tag{4.2.13}$$

We enforce normalization and the point-wise constraint $\Delta w = Ra\partial_x\theta$ via the Lagrange
multipliers $\lambda/2$ and $\gamma(x,z)/Ra$ respectively. Thus

$$\lambda_{min} = \inf_\theta \mathscr{F}$$

51

where $\mathscr{F} = \int \left[ |\nabla \theta|^2 + \theta w \tau' - \frac{\lambda}{2} \left(\theta^2 - L\right) + \frac{\gamma(x,z)}{Ra}\left(\Delta w - Ra\frac{\partial^2 \theta}{\partial x^2}\right) \right] \mathrm{d}x\mathrm{d}z. \qquad (4.2.14)$

Integrating by parts and using the fact that $\gamma(x,z)$ is also periodic in $x$, we arrive at an alternative form of $\mathscr{F}$ :

$$\mathscr{F} = \int \left[ -\theta\Delta\theta + \theta w \tau' - \frac{\lambda}{2}\left(\theta^2 - L\right) + \frac{\gamma(x,z)}{Ra}\Delta w - \theta\frac{\partial^2 \gamma}{\partial x^2} \right] \mathrm{d}x\mathrm{d}z. \qquad (4.2.15)$$

Given these two forms, a straightforward calculation shows that the Euler-Lagrange equations

$$\frac{\delta\mathscr{F}}{\delta\theta} = 0 \qquad (4.2.16)$$

$$\frac{\delta\mathscr{F}}{\delta w} = 0 \qquad (4.2.17)$$

$$\frac{\delta\mathscr{F}}{\delta\gamma} = 0 \qquad (4.2.18)$$

reduce to:

$$-2\Delta\theta + w\tau' - \frac{\partial^2\gamma}{\partial x^2} = \lambda\theta \qquad (4.2.19)$$

$$Ra\theta\tau' + \Delta\gamma = 0 \qquad (4.2.20)$$

$$\Delta w - Ra\frac{\partial^2\theta}{\partial x^2} = 0. \qquad (4.2.21)$$

Now, we show that the local extrema of $\mathscr{H}_1\{\theta\}$ are precisely the values of $\lambda$ that solve 4.2.19-4.2.21. Multiplying (4.2.19) by $\theta$ and integrating, we obtain

$$\mathscr{H}_1\{\theta\} - \int \left[ \theta\Delta\theta + \theta\frac{\partial^2\gamma}{\partial x^2} \right] \mathrm{d}x\mathrm{d}z = \lambda\int \theta^2\mathrm{d}x\mathrm{d}z = \lambda. \qquad (4.2.22)$$

52

However, equations (4.2.20) and (4.2.21) and the application of integration by parts show that

$$\int \left[\theta\Delta\theta + \theta\frac{\partial^2\gamma}{\partial x^2}\right] \mathrm{d}x\mathrm{d}z = \int \left[\theta\Delta\theta + \gamma\frac{\partial^2\theta}{\partial x^2}\right] \mathrm{d}x\mathrm{d}z \tag{4.2.23}$$

$$= \int \left[\theta\Delta\theta + \frac{\gamma}{Ra}\Delta w\right] \mathrm{d}x\mathrm{d}z \tag{4.2.24}$$

$$= \int \left[\theta\Delta\theta + \frac{w}{Ra}\Delta\gamma\right] \mathrm{d}x\mathrm{d}z \tag{4.2.25}$$

$$= \int \left[\theta\Delta\theta - w\theta\tau'\right] \mathrm{d}x\mathrm{d}z \tag{4.2.26}$$

$$= -\mathscr{H}_1\{\theta\}. \tag{4.2.27}$$

Therefore, $\mathscr{H}_1\{\theta\} = \lambda/2$ if $\theta$ and $\lambda$ solve $(4.2.19) - (4.2.21)$.

## 4.2.5 The adapted basis

It is a straightforward exercise to show that the functions $\theta(x, z)$ solving $(4.2.19 - 4.2.21)$ are in fact eigenfunctions of the symmetric operator $\mathsf{S}$. One can see this by noting that

$$\delta \int \theta\mathsf{S}\{\theta\} - \lambda |\theta|^2 = 0 \implies \mathsf{S}\{\theta\} = \lambda\theta. \tag{4.2.28}$$

This is the essence of our variational problem while the additional Lagrange multiplier serves merely as a means to implicitly enforce the momentum equation.

To compute the eigenfunctions, we take the Fourier transform of the equations in $x$ to arrive

at a $z$-dependent equation for each horizontal wavenumber $n$. Substitute:

$$\begin{bmatrix} \theta(x,z) \\ w(x,z) \\ \gamma(x,z) \end{bmatrix} \longrightarrow \begin{bmatrix} \Theta_{mn}(z)e^{inkx} \\ W_{mn}(z)e^{inkx} \\ \Gamma_{mn}(z)e^{inkx} \end{bmatrix}, \quad \|\Theta_{mn}\| = 1. \tag{4.2.29}$$

where $k = 2\pi/L$. The equations therefore become

$$(-2)\left[D^2 - (nk)^2\right]\Theta_{mn} + \tau' W_{mn} + (nk)^2\Gamma_{mn} = \lambda_{mn}\Theta_{mn} \tag{4.2.30}$$

$$\left[D^2 - (nk)^2\right]W_{mn} + (nk)^2 Ra\Theta_{mn} = 0, \tag{4.2.31}$$

$$\left[D^2 - (nk)^2\right]\Gamma_{mn} + Ra\tau'\Theta_{mn} = 0, \tag{4.2.32}$$

where the discrete set of vertical eigenfunctions for a given $n$ is indexed by the vertical "wave number" $m$. We note that for $n = 0$, The equations may be solved exactly:

$$W_{m0}(z) = 0, \quad \Theta_{m0}(z) = \sqrt{2}\sin((m+1)\pi z), \quad \lambda_{m0} = 2\left(m+1)\right)^2\pi^2 \tag{4.2.33}$$

## 4.2.6 Derivation of the ODES

The derivation of the Galerkin approximations of the equation of motion are similar to those we demonstrated in the case of the Fourier basis. The fields are expanded in the obtained eigenbasis with time-dependent amplitudes:

$$\theta(x,z) = \sum_{\beta=-\infty}^{\infty}\sum_{\alpha=0}^{\infty} a_{\alpha\beta}(t)\Theta_{\alpha\beta}(z)e^{i\beta kx}, \tag{4.2.34}$$

$$w(x,z) = \sum_{\beta=-\infty}^{\infty}\sum_{\alpha=0}^{\infty} b_{\alpha\beta}(t)W_{\alpha\beta}(z)e^{i\beta kx}, \tag{4.2.35}$$

$$u(x,z) = \sum_{\beta=-\infty}^{\infty}\sum_{\alpha=0}^{\infty} c_{\alpha\beta}(t)U_{\alpha\beta}(z)e^{i\beta kx}, \tag{4.2.36}$$

$$\gamma(x,z) = \sum_{\beta=-\infty}^{\infty}\sum_{\alpha=0}^{\infty} d_{\alpha\beta}(t)\Gamma_{\alpha\beta}(z)e^{i\beta kx}, \tag{4.2.37}$$

Here, the main difference is that the linear operator in 4.2.1 needs to be decomposed into its symmetric and antisymmetric parts.

Denoting by $w\{\theta\}$ the self-adjoint[3] linear operator such that $w\{\theta\}(x,z,t) = w(x,z,t)$, we compute $\mathsf{S}$ and $\mathsf{A}$ as follows. For any function $f$ satisfying the same boundary conditions as $\theta$,

$$(f, \mathscr{L}\{\theta\}) = \int \left[ -f\Delta\theta + f\tau'w\{\theta\}\right]\mathrm{d}x\mathrm{d}z \tag{4.2.38}$$

$$= \int \left[ -\Delta f\,\theta + w\{f\tau'\}\,\theta\right]\mathrm{d}x\mathrm{d}z. \tag{4.2.39}$$

Thus,

$$\mathscr{L}^{\dagger}\{\theta\} = -\Delta\theta + w\{\tau'\theta\} \tag{4.2.40}$$

[3]One can see this from the self-adjointness of the inverse Laplacian operator and the momentum equation.

which yields

$$\mathsf{S}\{\theta\} = \frac{1}{2}\left(\mathscr{L}\{\theta\} + \mathscr{L}^{\dagger}\{\theta\}\right) = -\Delta\theta + \frac{1}{2}\left(\tau' w\{\theta\} + w\{\tau'\theta\}\right) \qquad (4.2.41)$$

$$\mathsf{A}\{\theta\} = \frac{1}{2}\left(\mathscr{L}\{\theta\} - \mathscr{L}^{\dagger}\{\theta\}\right) = \frac{1}{2}\left(\tau' w\{\theta\} - w\{\tau'\theta\}\right). \qquad (4.2.42)$$

Equation 4.2.1 then becomes

$$\partial_t\theta + \boldsymbol{u}\cdot\boldsymbol{\nabla}\theta = \tau'' - \mathsf{S}\{\theta\} - \mathsf{A}\{\theta\}. \qquad (4.2.43)$$

Once more, we denote the eigenfunctions by $\phi_{mn}(x,z) = \Theta_{mn}(z)e^{inkx}$ and project the time-dependent equation onto each $\phi_{mn}$. By the momentum equation and incompressibility respectively, we have the two constraints $b_{qj}(t) = a_{qj}(t)$ and $c_{mn}(t)U_{mn}(z) = (i/nk)b_{mn}(t)W_{mn}(z)$.

$$(\phi_{mn}, \partial_t\theta) = \dot{a}_{mn} \qquad (4.2.44)$$

$$(\phi_{mn}, \mathsf{S}\{\theta\}) = -\frac{\lambda_{mn}}{2}a_{mn} \qquad (4.2.45)$$

$$(\phi_{mn}, \tau'') = \int \tau''(z)\Theta_{mn}(z)\delta_{n,0}\mathrm{d}z \qquad (4.2.46)$$

$$(\phi_{mn}, \boldsymbol{u}\cdot\boldsymbol{\nabla}\theta) = \sum_{\beta,q=-\infty}^{\infty}\sum_{\alpha,p=0}^{\infty}(ik\beta)c_{p,q}a_{\alpha,\beta}\int\Theta_{m,n}U_{p,q}\Theta_{\alpha,\beta}\delta(\beta+q-n)\mathrm{d}z$$
$$+ \sum_{\beta,q=-\infty}^{\infty}\sum_{\alpha,p=0}^{\infty}b_{p,q}a_{\alpha,\beta}\int\Theta_{m,n}W_{p,q}D\Theta_{\alpha,\beta}\delta(\beta+q-n)\mathrm{d}z \quad (4.2.47)$$

After simplification, the application of modal constraints, and finally renaming the dummy

variables in the sums, the quadratic term simplifies to

$$(\phi_{mn}, \boldsymbol{u} \cdot \boldsymbol{\nabla}\theta) \;=\; \sum_{\substack{j \neq n \\ j=-\infty}}^{\infty} \sum_{p=0}^{\infty} \sum_{q=-\infty}^{\infty} a_{qj} q_{p(n-j)} \int \left[ \left( \frac{j}{n-j} \right) DW_{p(n-j)} \Theta_{qj} - W_{p(n-j)} D\Theta_{qj} \right] \Theta_{mn} \mathrm{d}z$$

$$\text{if } n \neq 0 \tag{4.2.48}$$

and

$$(\phi_{mn}, \boldsymbol{u} \cdot \boldsymbol{\nabla}\theta) \;=\; \sum_{\substack{j \neq 0 \\ j=-\infty}}^{\infty} \sum_{p=0}^{\infty} \sum_{q=0}^{\infty} a_{qj} a_{p(n-j)} \sqrt{2}(m+1)\pi \int_0^1 W_{pj} \Theta_{qj} \cos((m+1)\pi z)\, \mathrm{d}z$$

$$\text{if } n = 0. \tag{4.2.49}$$

Finally, we project the antisymmetric term onto $\phi_{mn}$. From $(4.2.20)$ we find that $\Delta \left( \frac{\partial^2 \gamma}{\partial x^2} \right) = Ra \frac{\partial^2}{\partial x^2} \left( \tau' \theta \right)$ which implies that $w\{\tau'\theta\} = \frac{\partial^2 \gamma}{\partial x^2}$. Combining this with $(4.2.42)$, we deduce

$$(\phi_{mn}, \mathsf{A}\{\theta\}) = -\frac{1}{2} \sum_{p=0}^{\infty} (nk)^2 d_{pn} \int_0^1 \Gamma_{pn} \Theta_{mn} \mathrm{d}z + \frac{1}{2} \sum_{p=0}^{\infty} b_{pn} \int_0^1 \tau' W_{pn} \Theta_{mn} \mathrm{d}z. \tag{4.2.50}$$

Using $(4.2.32)$, integration by parts, and then $(4.2.31)$, we find

$$\frac{1}{2} \sum_{p=0}^{\infty} b_{pn} \int_0^1 \tau' W_{pn} \Theta_{mn} \mathrm{d}z \;=\; -\frac{1}{2} \sum_{p=0}^{\infty} b_{pn} \frac{1}{Ra} \int_0^1 \left[ W_{pn} D^2 \Gamma_{mn} - (nk)^2 W_{pn} \Gamma_{mn} \right] \mathrm{d}z$$

$$= -\frac{1}{2Ra} \sum_p b_{pn} \int_0^1 \left[ D^2 - (nk)^2 \right] W_{pn} \Gamma_{mn} \mathrm{d}z$$

$$= \frac{1}{2} \sum_p b_{pn} \int_0^1 (nk)^2 \Theta_{pn} \Gamma_{mn} \mathrm{d}z. \tag{4.2.51}$$

57

Therefore,

$$(\phi_{mn}, \mathsf{A}\{\theta\}) = \frac{1}{2}\sum_p a_{mn}(nk)^2 \int_0^1 [\Theta_{pn}\Gamma_{mn} - \Theta_{mn}\Gamma_{pn}]\,\mathrm{d}z. \tag{4.2.52}$$

Combining all terms, the final form of the ordinary differential equations will be

$$\dot{a}_{mn} = \mu_{mn}a_{mn} + \sum_{p=0}^{\infty}\mu_{mn}^p a_{pn} + \sum_{j=-\infty}^{\infty}\sum_{p=0}^{\infty}\sum_{q=0}^{\infty}\Lambda_{mn}^{jpq}a_{qj}a_{p(n-j)} \quad \text{for } n \neq 0 \tag{4.2.53}$$

$$\dot{a}_{m0} = \mu_{m0}a_{m0} + \sum_{j=-\infty}^{\infty}\sum_{p=0}^{\infty}\sum_{q=0}^{\infty}\Lambda_{m0}^{jpq}a_{qj}a_{pj}^* + f_m \tag{4.2.54}$$

where

$$\mu_{mn}^p = \left(\frac{n^2 k^2}{2}\right)\int_0^1 [\Theta_{mn}\Gamma_{pn} - \Theta_{pn}\Gamma_{mn}]\,\mathrm{d}z, \tag{4.2.55}$$

$$\mu_{m0} = -(m+1)^2\pi^2, \tag{4.2.56}$$

$$\mu_{mn} = -\frac{\lambda_{mn}}{2} \quad \text{for } n \neq 0. \tag{4.2.57}$$

$$\Lambda_{mn}^{jpq} = \int_0^1 \Theta_{mn}\left[\left(\frac{j}{j-n}\right)DW_{p(n-j)}\Theta_{qj} - W_{p(n-j)}D\Theta_{qj}\right]\,\mathrm{d}z, \tag{4.2.58}$$

$$\Lambda_{m0}^{jpq} = \sqrt{2}(m+1)\pi\int_0^1 W_{pj}\Theta_{qj}\cos((m+1)\pi z)\,\mathrm{d}z, \tag{4.2.59}$$

$$f_m = \sqrt{2}(m+1)\pi\left[1 - (m+1)\pi\int_0^1 \tau(z)\sin((m+1)\pi z)\,\mathrm{d}z\right], \tag{4.2.60}$$

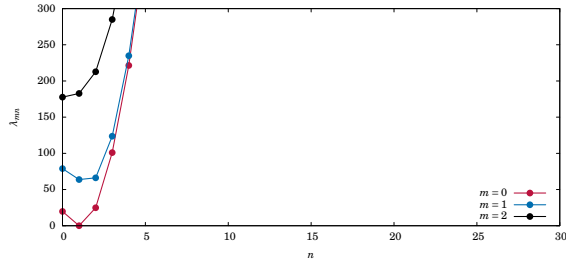## 4.2.7 Numerical computation of the spectrum

It is possible to solve the constrained variational problem (4.2.12) by brute force optimization using standard numerical optimization software packages as in [39, 40]. There, we computed the spectra and the associated upper bounds for up to $Ra = 2102$. However for $Ra$ greater than a few thousand, the method proves excessively resource intensive and lacking in robust-

ness. Numerical instabilities prohibit accurate computation of the optimal eigenfunctions. In [41] a new strategy was devised for efficiently and accurately solving the variational problem which extends the results to $Ra \approx 2.65 \times 10^4$. This method consists of two steps: first an Euler-Lagrange problem is directly formulated for the minimization of $\int_0^1 \tau'^2(z) \, dz$ subject to the spectral constraint. The resulting Euler- Lagrange equations are augmented by a time-derivative term to yield time evolution equations with steady states equal to the desired solutions. A time-marching method is then used to find the critical modes (those with $\lambda = 0$). Good approximations for the critical modes are quickly found. However, convergence to the full solution may be forbiddingly slow. Therefore, in the second step, initial guesses for a Newton-Kantorovich (NK) iterative method are constructed using only the critical modes found in the first step and the full solution is then computed using the NK method.

The spectrum thus computed possesses properties pertinent to low-dimensional modeling. To demonstrate this, we present the spectrum and eigenfunctions computed for several values of $Ra$ in fig 4.2.1.

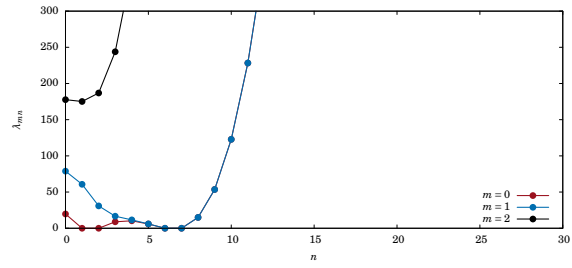Several features of the optimal spectra are worth noting:

1. The spectra indicate the existence of "marginally stable" modes, i.e., modes with $\lambda = 0$. In dynamical terms, these modes will have zero linear damping and will be driven entirely by nonlinear interaction( and constant forcing represented by the $F_m$ terms) whereas all other modes will be damped linearly. As $Ra$ increases, more and more marginally stable modes are introduced and/or they shift toward higher wave numbers, in line with the notion that the dominant length scales decrease as $Ra$ increases.

2. There is a distinct separation between the first two branches of the spectrum ($m = 0, 1$) and the rest. The first two branches tend to be drawn toward zero to to make
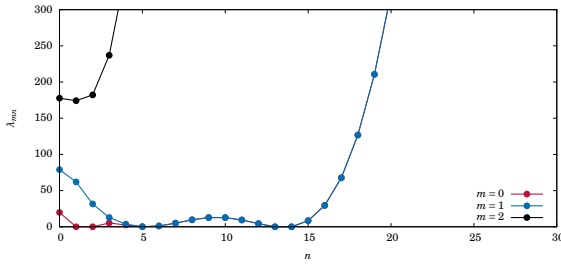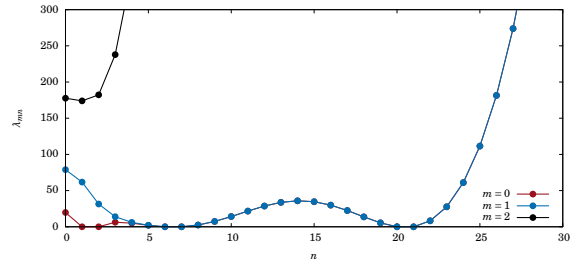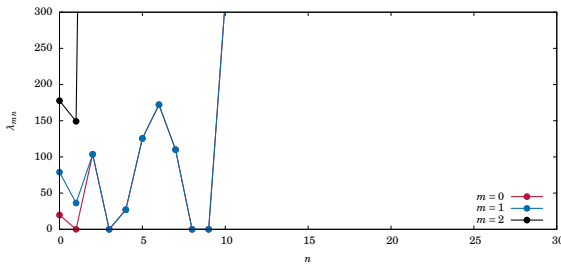
59

(a) $Ra = 100$, $L = 2$

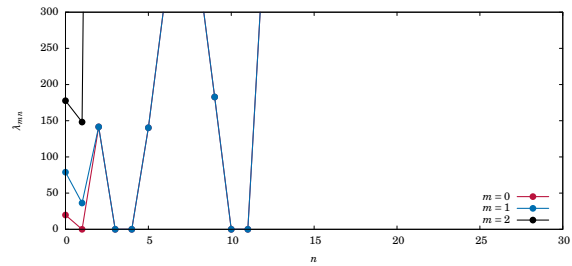(b) $Ra = 700$, $L = 2$

(c) $Ra = 900$, $L = 2$

(d) $Ra = 2000$, $L = 2$
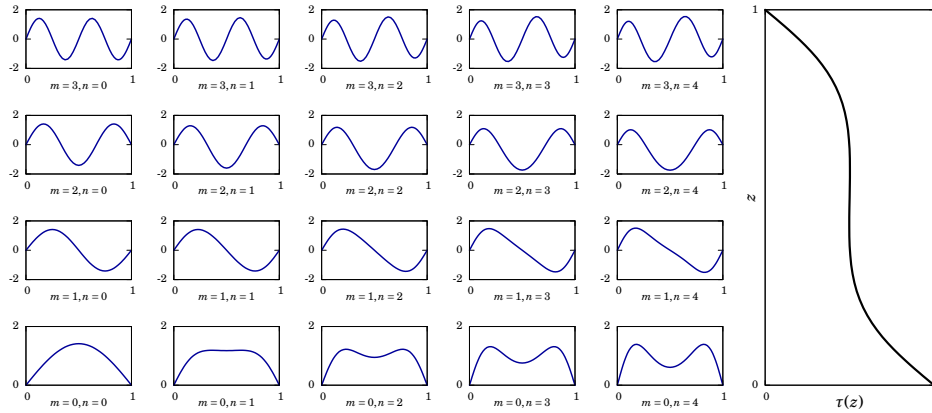
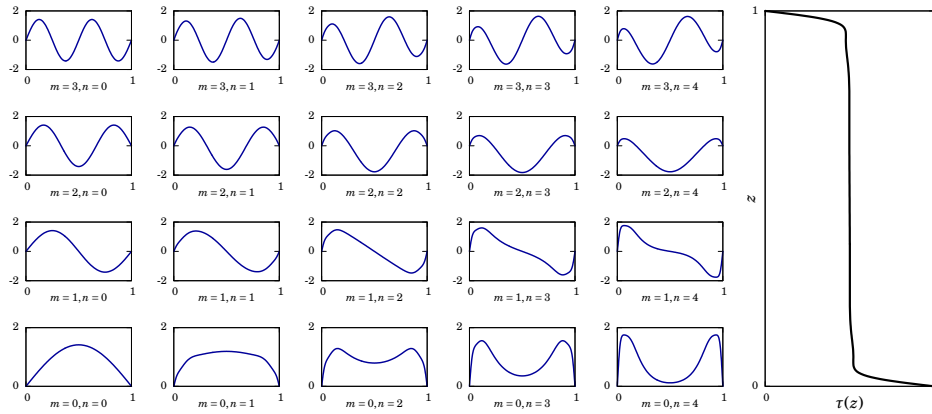(e) $Ra = 3000$, $L = 2$

(f) $Ra = 7000$, $0.3691$

(g) $Ra = 10000$, $L = 0.3156$

Figure 4.2.1: Numerically computed eigenvalues $\lambda_{mn}$ at various $Ra$ values. Notice the increasing number of marginally stable modes ($\lambda_{mn} = 0$) and their shift toward higher wave numbers as $Ra$ increases.
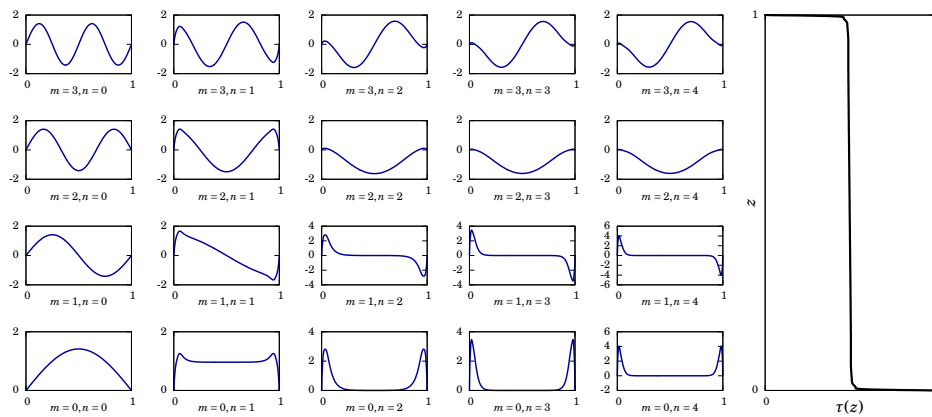
(a) $Ra = 100$

(b) $Ra = 1000$

(c) $Ra = 10000$

Figure 4.2.2: Some temperature eigenfunctions $\Theta_{mn}(z)$ and the optimal background profile $\tau(z)$ for $Ra = 100, 1000, 10000$. As $Ra$ increases, $\tau(z)$ mimics the horizontally averaged temperature with increasingly thin boundary layers. The temperature eigenvalues also tend to concentrate increasingly toward $z = 0, 1$. 61

"touchdowns" a number of times before growing unboundedly. Remarkably, the second branch eventually mimics the first, producing marginally stable modes at the same wavelengths. The alternating parity of $\Theta_{mn}$ about $z = 1/2$ with increasing $m$ (evident in (4.2.2$c$) for instance) means that using modes from the first two branches alone, one can produce independent dynamics at the two boundary layers.

3. The optimal background profile $\tau(z)$ resembles the horizontally averaged temperature in that the sharpest gradients occur near $z = 0, 1$ and there is little variation in the bulk. Furthermore, with increasing $Ra$, we observe a thinning of the boundary layers in $\tau$, reminiscent of the actual (time-averaged) temperature empirically and numerically observed . Thus, in our Nonlinear Galerkin models, the boundary layer is more or less resolved a priori by the choice of $\tau$, removing the burden from the dynamics.

We also note that new upper bounds for $Nu$ may be obtained as a bi-product of the numerical computation of the optimal background profile. Once the optimal background profile $\tau(z)$ is computed, using (4.2.5) an upper bound (at $Ra/2$) may be found for any $a \in (0, 1)$. Finally, the optimal upper bound is computed numerically by varying the parameter $a$. Figure 4.2.3 shows the new upper bounds[4] alongside previously obtained analytical and numerical bounds, and results from direct numerical simulations.

## 4.3   Discussion

A number of qualitative and quantitative distinctions can be predicted between the FG and NG models. For instance, we can derive a simple absolute lower bound on the size of any

---

[4]The numerical computation of the optimal bounds is due to Baole Wen and Gregory Chini, University of New Hampshire.
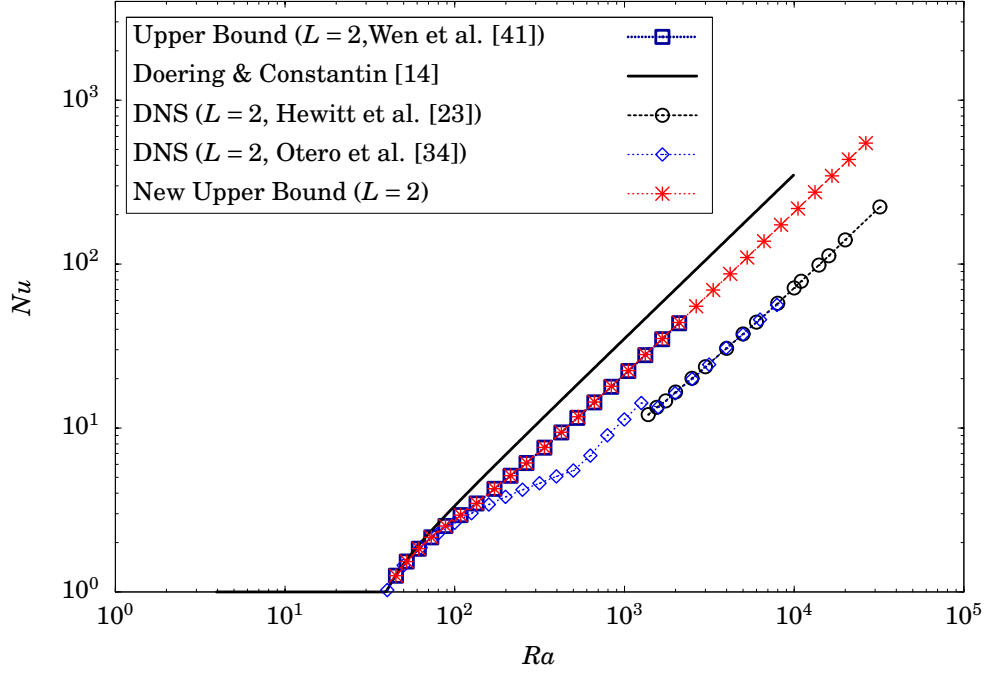
Figure 4.2.3: Previously obtained analytical and numerical upper bounds, compared with DNS results and the new upper bounds.

NG model as follows. The Galerkin-truncated equations of motion are

$$\partial_t \theta_N + P_N \left[ \boldsymbol{u}_N \cdot \boldsymbol{\nabla} \theta_N \right] \;=\; \Delta \theta_N + w_N \qquad (4.3.1)$$

$$\Delta w_N \;=\; Ra \frac{\partial^2 \theta_N}{\partial x^2} \qquad (4.3.2)$$

$$\boldsymbol{\nabla} \cdot \boldsymbol{u}_N \;=\; 0. \qquad (4.3.3)$$

where $\theta_N$ is the $N$-th Galerkin approximation of $\theta$ and $w_N$ and $u_N$ solve (4.3.2) and (4.3.3). Note that even though $\boldsymbol{u}_N \neq \boldsymbol{u}$, it is still incompressible by design. Thus, if we multiply (4.3.1) by $\theta_N$ and integrate, the cubic term vanishes just as in (4.2.11) with the full fields. The expression for the rate of energy change then contains quadratic terms only.

Therefore, we expect all cubic terms to vanish also when we repeat the computation using the Fourier-Galerkin ODEs (4.1.29). The only remaining terms are those originating from the linear term of the time evolution equations:

$$\frac{1}{2}\frac{d}{dt}\|\theta_N\|_2^2 = \frac{1}{2}\frac{d}{dt}\sum|a_{mn}|^2 \tag{4.3.4}$$

$$= \frac{1}{2}\sum[\dot{a}_{mn}a_{mn}^* + \dot{a}_{mn}^*a_{mn}] \tag{4.3.5}$$

$$= \operatorname{Re}\{a_{mn}^*\dot{a}_{mn}\}. \tag{4.3.6}$$

Thus

$$\frac{1}{2}\frac{d}{dt}\|\theta_N\|_2^2 = \sum \mu_{mn}|a_{mn}|^2. \tag{4.3.7}$$

This implies that energy will diverge unless at least one stable mode is included in the basis. In other words, the number of unstable modes is a lower bound on the size of any working model. Of course, the inclusion of at least one stable mode is only a necessary condition and does not guarantee energy-stability, let alone any measure of convergence. In order to estimate this lower bound as a function of $Ra$, we count the modes with $\mu_{mn} > 0$ (4.1.27):

$$\frac{Ra\,m^2}{m^2 + n^2} > \pi^2(m^2 + n^2) \quad \text{with} \quad m, n > 0 \tag{4.3.8}$$

leads to $-m^2 + \sqrt{\frac{Ra}{\pi^2}}m > n^2$ and thus, to leading order, the number of unstable modes is

$$\int_0^{\sqrt{\frac{Ra}{\pi^2}}} \left(-m^2 + \sqrt{\frac{Ra}{\pi^2}}m\right)^{1/2} dm = \frac{Ra}{8\pi} \tag{4.3.9}$$

as $Ra \to \infty$. At $Ra = 1000$ for instance, this implies that any model with less than $\sim 40$ modes will necessarily diverge.

Next chapter will be dedicated to numerical studies of dynamical systems of various truncations, using both FG and NG.

# CHAPTER V

# Low-dimensional models

In the last chapter, we derived two classes of Galerkin spectral methods for the solution of porous-medium convection: the standard Fourier-Galerkin methods obtained from the Galerkin projection of the PDEs onto the Fourier basis, and the new "Nonlinear-Galerkin" methods that exploit the "adapted basis" and the optimal background profile produced by our proposed variational scheme.

In this chapter, we will present and discuss the results of a substantial set of simulations of both methods at various Rayleigh numbers. The goal is to understand the relative strengths of the two methods in different dynamical regimes. In particular, we investigate the manner in which the solutions of the two methods "converge" as more and more modes are included. For added perspective, we also compare some results with those of pseudo-spectral collocation methods as described in Chapter III.

## 5.1 Finite truncations

A well-defined finite-dimensional dynamical model (Galerkin method) has not been constructed until we have truncated the infinity of ODEs to include only a finite number of selected modes. For extreme truncations where only a few modes are kept, one might be able to select the modes based on heuristic physical interpretations of the role of those modes in the dynamics. However, in general one needs a consistent rule that yields an increasing sequence of sets of modes to be used for models of various "sizes". Questions of convergence may then be studied in a systematic fashion.

Given the ODEs computed in the previous chapter, the most immediate measure of the comparative dynamical vigor of different modes is the linear growth coefficient of each mode. In fact, linear analysis shows that at the onset of convection, the linear coefficients are solely responsible for determining the unstable modes and their relative growth rates. However, away from the conduction solution, the nonlinear interactions dominate and one can not rely on linear arguments any more. The other coefficients in the equations define linear or quadratic couplings between the modes and therefore, from a pragmatic point of view, do not contribute to an ordering of the modes in any trivial way. On the other hand, the linear coefficients quantify the level of self-inhibition (when negative) or linear growth (when positive) and thus provide a clear albeit incomplete measure of dynamical relevance. We rely solely on the linear coefficients for defining our truncations.

Specifically, let the tuples $\{(m_1, n_1), (m_2, n_2), \cdots\}$ index the modes such that

$$\mu_{m_1 n_1} \geq \mu_{m_2 n_2} \geq \mu_{m_3 n_3} \geq \cdots . \tag{5.1.1}$$

67

In our notation, for the NG model, $m = 0, 1, 2, \cdots$ and $n = 0, 1, 2, \cdots$ whereas in the FG model, $m = 1, 2, 3, \cdots$. A truncation of size $N$ is defined as the choice of the modes

$$\{(m_1, n_1), (m_2, n_2), \cdots, (m_N, n_N)\}. \tag{5.1.2}$$

Initially, we introduce a slight modification to this scheme. Since computationally the modes themselves are commonly computed on finite "rectangles" in the spectral space, i.e., $(m, n) \in \{M_{min}, \cdots, M_{max}\} \times \{N_{min}, \cdots, N_{max}\}$, and since additionally this is the only option in the case of pseudo-spectral methods, unless otherwise indicated, a truncation of size $N$ consists of the $N$ modes with the largest linear coefficients within a rectangle of predefined size in the spectral space. For instance, Figure 5.1.1 shows how this truncation is applied to the Fourier-Galerkin models.

## 5.2   Numerical methods

The implementation of both the NG and the FG models requires the numerical computation of all the coefficients $(4.2.55 - 4.2.60)$ and $(4.1.23 - 4.1.27)$ pertaining to the truncation first. The coefficients derived for the FG model are algebraic functions of various modal indices and thus efficiently computable. On the other hand, those of the NG model, require differentiation and integration of various combinations of the numerically computed eigenfunctions on the Lobatto grid. This task may be performed with high precision, but is relatively expensive. It should be noted that overall, the numerical computation of the NG equations takes tens of times longer than that of the FG equations. We have taken

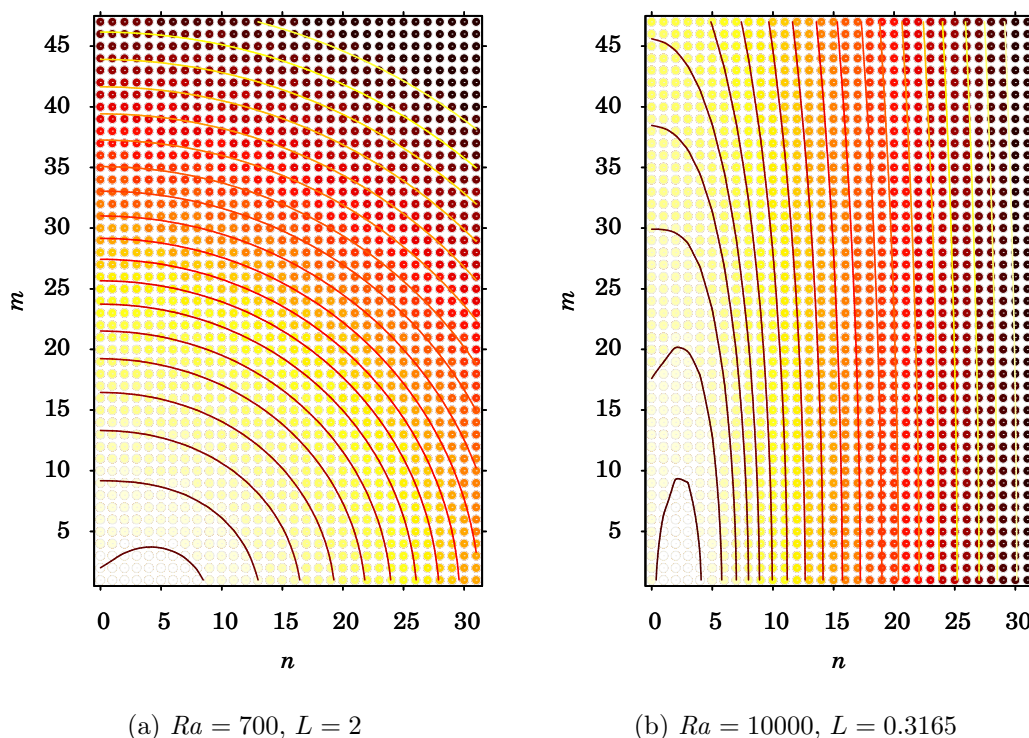(a) $Ra = 700$, $L = 2$      (b) $Ra = 10000$, $L = 0.3165$

Figure 5.1.1: An illustration of the truncation scheme for the Fourier-Galerkin method. The color encodes the linear coefficient of the mode, with lighter colors indicating larger values and thus higher precedence. Each contour delineates the set of modes used for a truncation. The truncations depicted use increasing percentages of the total number of modes, equally spaced from 2 percent to 100 percent.

advantage of parallelized algorithms running simultaneously on multiple cores [5] to perform these computations. Once the equations are computed for the largest desired truncation, those of lower truncations may be extracted as appropriately chosen subsets.

We performed the time integration of the equations using a fourth order Runge-Kutta method with adaptive time-stepping [42]. For each batch of simulations at a given $Ra$, a single run at a relatively low truncation was performed to ensure the desired type of solution was

---

[5] Using the FLUX computer cluster at the University of Michigan.

produced, and the result was used as the initial condition for all other truncations. This involves using the old modal amplitudes in the new truncation and assigning zero to the new modes. At high $Ra$, the presence of the zero amplitudes in the initial conditions can lead to numerical instability unless the time steps are initially chosen to be several orders of magnitude smaller than otherwise necessary so that the new modes are populated gradually. Hence the use of adaptive time-stepping.

The Nusselt number is computed using the bulk formula (2.2.23) unless otherwise stated.

At low to medium $Ra$, i.e., $Ra = 100, 700, 900$, the solution we seek to simulate is the single pair of rolls (steady for $Ra = 100$ and with boundary layer plumes at $Ra = 700, 900$). At $Ra = 100$, this solution is centro-symmetric modulo a horizontal translation, i.e., if centered in the box, each half-box of size $1 \times 1$ appears symmetric about $x = z = \frac{1}{2}$. At $Ra = 700, 900$, there exists a family of solutions, all organized about the single pair of rolls with boundary layer plumes, but not all need be centro-symmetric. However, our results show no difference in $Nu$ between the two. A restriction to the subspace of centro-symmetric solutions is thus physically justified and computationally advantageous: one needs to model only half of the modes in the centro-symmetric subspace. For more details on centro-symmetry, see Appendix B.


## 5.3   Direct numerical simulations

Before we begin the simulation of the Galerkin models, we first examine the convergence of the standard pseudo-spectral collocation method described in Chapter III. Here, the goal is to solve the porous-medium convection problem at several different grid sizes, which we interpret in terms of the number of "modes" used, and observe the pattern of convergence

for the resulting $Nu$.


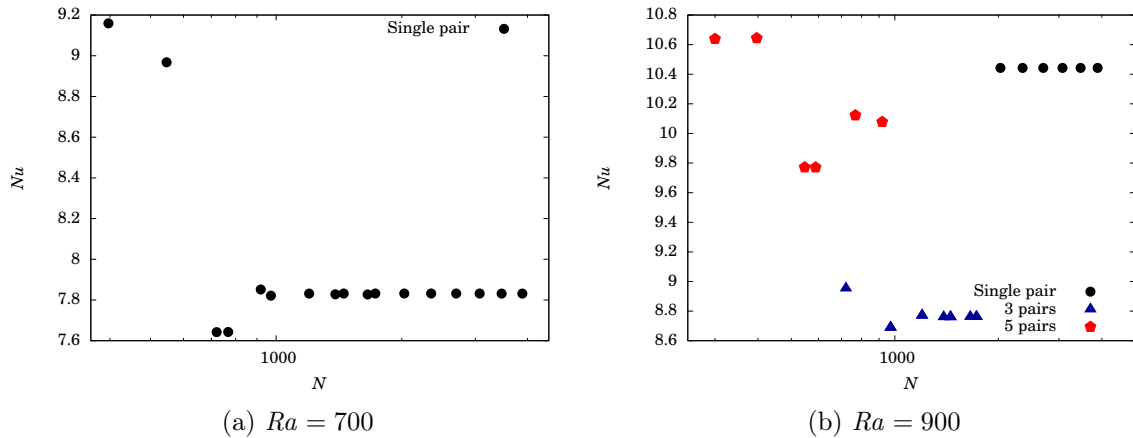
(a) $Ra = 700$

(b) $Ra = 900$

Figure 5.3.1: Pseudo-spectral simulations at $Ra = 700$ and $Ra = 900$ at different resolutions. The former appears to converge with about 1000 modes and the latter with $1000 \sim 2000$ modes.

We perform these simulations at $Ra = 700$ and $Ra = 900$. Following Graham and Steen [21], the vertical and horizontal resolutions of our box of $L = 2$ are $N_z$ and $\frac{8}{3}N_z$ respectively, where $22 \leq N_z \leq 72$. We restrict the solutions to the centro-symmetric subspace. Thus, by centro-symmetry and reality, the total number of degrees of freedom is $N = \frac{2}{3}N_z^2$. In both $Ra$ values, we find other solutions as well, including a steady three pairs of rolls. Especially at $Ra = 900$, it is not always easy to isolate the single-pair solution. Thus, given our data, we can determine the size of the smallest converging simulation only approximately. The results (Figure 5.3.1) indicate that at these $Ra$ values, the simulations converge with 1000-2000 degrees of freedom. As we will discuss later, this is far more than the minimum number of degrees of freedom necessary for the convergence of either Galerkin method, illustrating the fundamental disparity between the two types of spectral methods. Having established this, we henceforth focus on comparing the two different Galerkin methods.
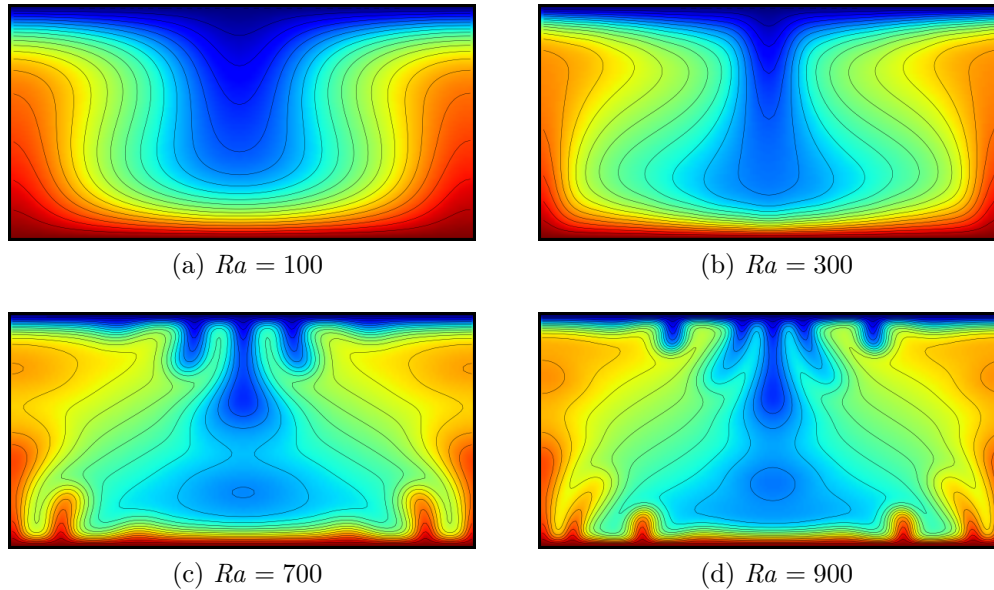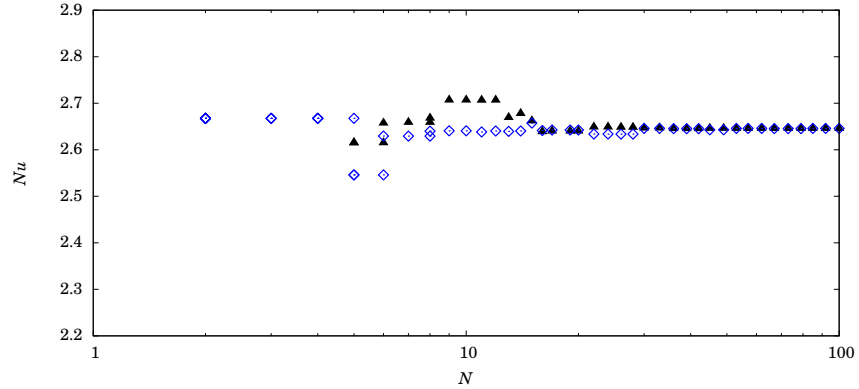
(a) $Ra = 100$      (b) $Ra = 300$

(c) $Ra = 700$      (d) $Ra = 900$

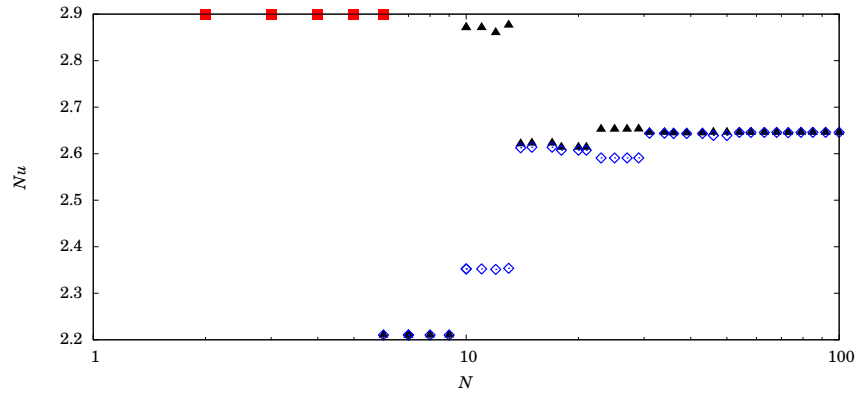Figure 5.3.2: Temperature fields from converged pseudo-spectral direct numerical simulation at low $Ra$.

## 5.4 Dynamical systems models

### 5.4.1 Low $Ra$ regime: $Ra = 100$

At $Ra = 100$, our simulations show a significant difference between the NG and FG models. As Figure 5.4.1 illustrates, The FG models diverge with $N < 6$, need at least 14 modes to yield a reasonably accurate $Nu$ and do not converge until $N \sim 30$. On the other hand, the NG models do not diverge even at $N = 2$, and at $N = 6$, then already produce a $Nu$ within a few percents of the converged value. They more or less converge at $N = 17$.
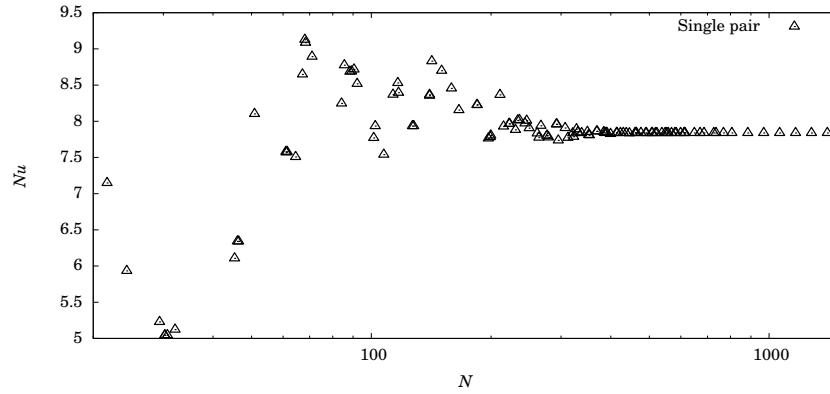
(a) Non-linear Galerkin models, $Ra = 100$.



(b) Fourier Galerkin models, $Ra = 100$.

Figure 5.4.1: Comparison of NG and FG models at various truncations. The triangles measure $Nu$ from conduction at the boundary layers whereas the diamonds use the bulk-averaged heat transport formula (2.2.23). The red squares indicate diverging models.
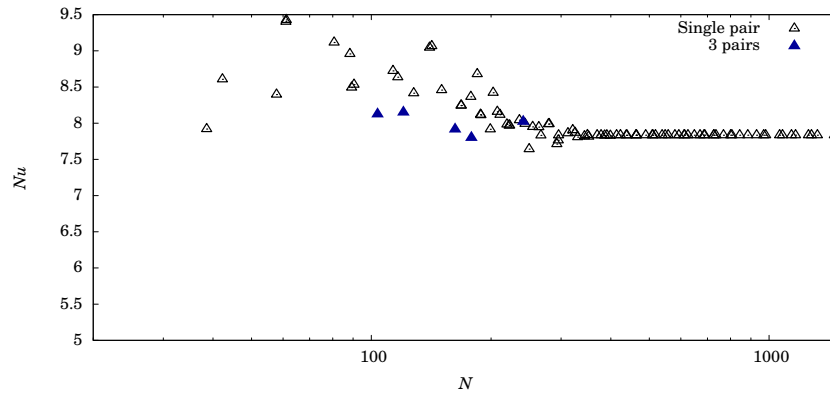
## 5.4.2 Medium $Ra$ regime: $Ra = 700$ and $Ra = 900$

At $Ra = 700$ and $Ra = 900$, we observe very little difference between the two methods. As figures 5.4.2 and 5.4.3 illustrate, at $Ra = 700$, both models converge at $N = 300 \sim 400$ whereas at $Ra = 900$, this occurs at $N = 500 \sim 600$. The only apparent difference is at extreme truncations ($N < 20 \sim 30$) where the FG model simply diverges. A major difficulty

in both models is isolating the single-pair solution at low $N$. The results presented in figures 5.4.2 and 5.4.3 have been chosen from hundreds of simulations most of which do not settle into the single-pair solution at low $N$ where meaningful differences between the two models are expected.
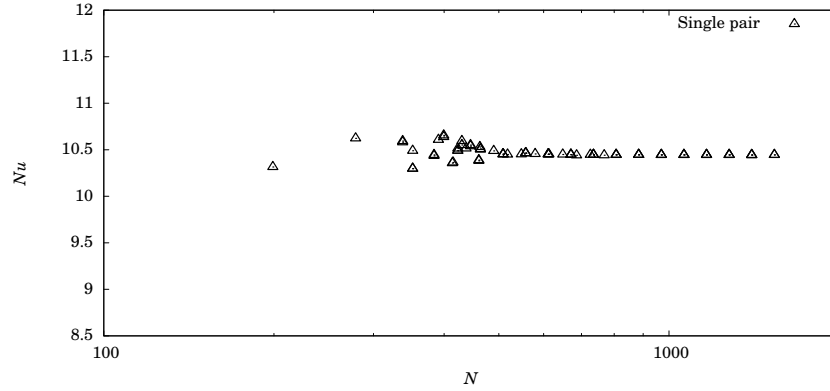


(a) Non-linear Galerkin models, $Ra = 700$.
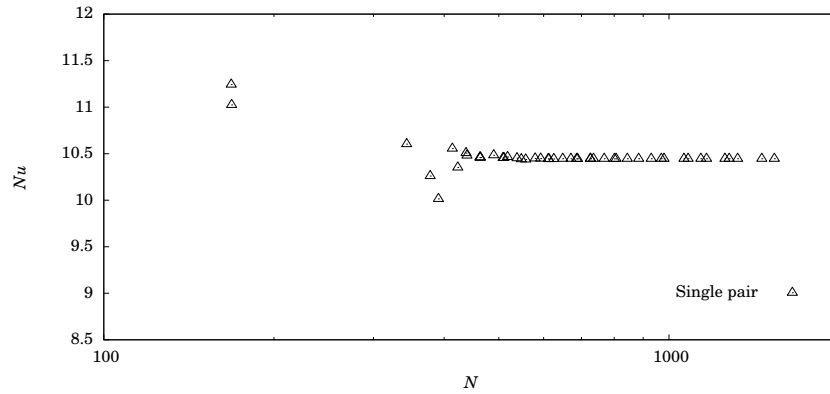


(b) Fourier Galerkin models, $Ra = 700$.

Figure 5.4.2: Comparison of NG and FG models at various truncations. $Ra = 700$.
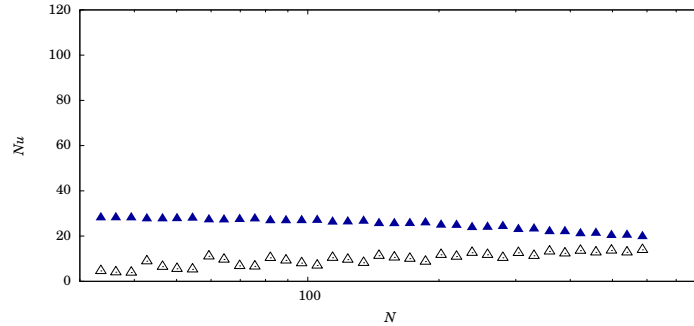
(a) Non-linear Galerkin models, $Ra = 900$.



(b) Fourier Galerkin models, $Ra = 900$.

Figure 5.4.3: Comparison of NG and FG models at various truncations. $Ra = 900$.

### 5.4.3  High $Ra$ regime: $Ra = 1500$, $L = 2$

As mentioned before, above $Ra \sim 1250$ the solution organized about a single pair of rolls becomes unstable, and the roll structure is replaced by multiple chaotic "mega-plumes". These solutions are not consistent with centro-symmetry and arise as the generic solution regardless of the initial conditions. Therefore, they are much easier to produce than the target solutions at $Ra = 700$ and $Ra = 900$. This allows us to study severe truncations and

compare the NG and FG methods.



(a) Non-linear Galerkin models, $Ra = 1500$.



(b) Fourier Galerkin models, $Ra = 1500$.

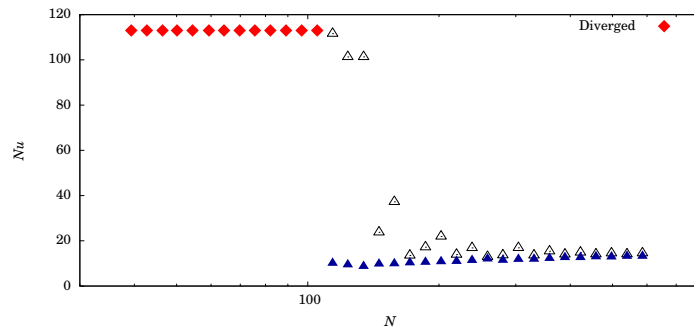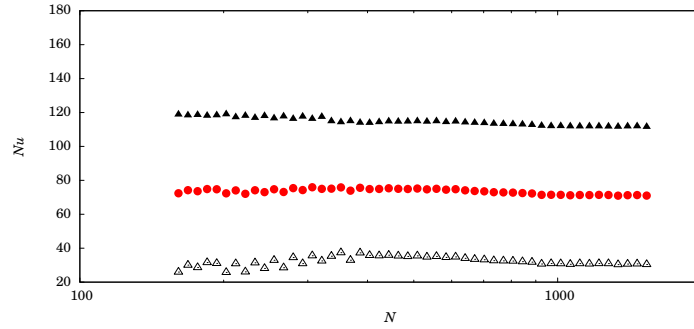Figure 5.4.4: Comparison of NG and FG models at various truncations. $Ra = 1500$, $L = 2$. Empty triangles represent $Nu$ computed using bulk-averaging whereas filled triangles indicated $Nu$ computed from conduction at boundary layers.

In this regime, the divergence of highly truncated FG models is rather nuanced. Beyond the diverging low-$N$ truncations, the two measurements of $Nu$—namely the one using the bulk average heat flux and the other using conductive heat flux at the boundaries— rapidly approach one another and continue to converge to the exact value. In contrast, the NG model produces finite Nusselt measurements using both definitions starting from very low $N$, but the two measurements approach one another at a much lower rate.

## 5.4.4 High $Ra$ regime, minimal flow unit:$Ra = 7 \times 10^3$ and $Ra = 10^4$



(a) Non-linear Galerkin models, $Ra = 7000$.



(b) Fourier Galerkin models, $Ra = 7000$.

Figure 5.4.5: Comparison of NG and FG models at various truncations. $Ra = 7000$, $L = 0.3691$. Empty triangles represent $Nu$ computed using bulk-averaging whereas filled triangles indicated $Nu$ computed from conduction at boundary layers. Solid circles represent the average of the two.
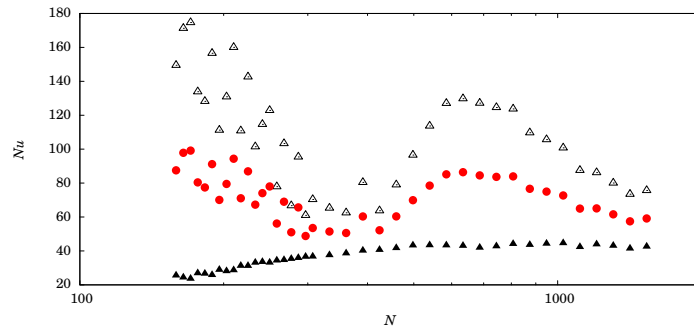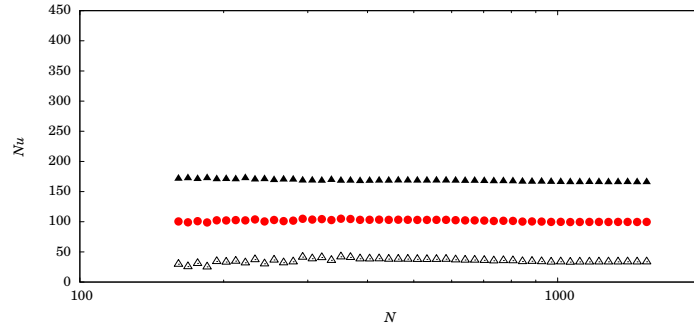
At $Ra = 7000$ and $Ra = 10000$ where the the flow is already turbulent, we exploit the emergence of the minimal flow unit to further reduce the size of our Galerkin models. Based on the conclusions of Chapter III, we compute the non-linear eigenfunctions and the resulting ODEs at $L = 4\pi Ra^{-2/5}$. Thus, the box sizes are $L = 0.3691$ and $L = 0.3156$ for $Ra = 7000$ and $Ra = 10000$ respectively. As illustrated by figure 4.2.1, reducing the box size leads to a more rapid (negative) growth of the linear modal coefficients and thus a smaller upper bound on the horizontal wave numbers included in truncations. Figure 5.4.7b also reflects
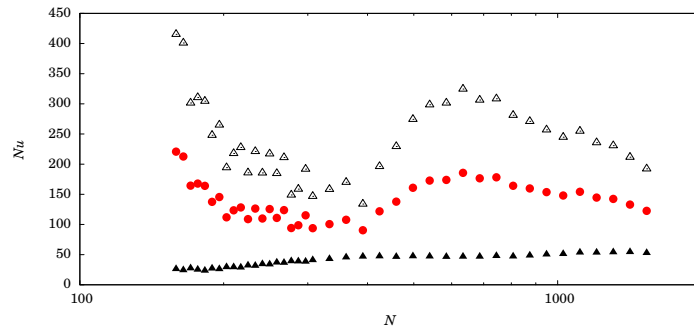
77

(a) Non-linear Galerkin models, $Ra = 10000$.



(b) Fourier Galerkin models, $Ra = 10000$.

Figure 5.4.6: Comparison of NG and FG models at various truncations. $Ra = 10000$, $L = 0.3156$. Empty triangles represent $Nu$ computed using bulk-averaging whereas filled triangles indicated $Nu$ computed from conduction at boundary layers. Solid circles represent the average of the two.

this fact.

Here, we are forced to forgo centro-symmetry for otherwise we are completely unable to isolate the desired solutions due to the strict symmetry imposed. Figure 5.4.6b summarizes our findings at this $Ra$. Surprisingly, in the NG case we find that $Nu$ remains more or less flat after a modest and brief initial move whereas for the FG case, we witness a sharp initial fall followed by a steep rise and finally a steady trend toward saturation.

(a) $Ra = 700$, $L = 2$

(b) $Ra = 10000$, $L = 0.3156$

(c) $Ra = 10000$, $L = 0.3156$, larger index ranges.

Figure 5.4.7: The truncation schemes for the non-linear Galerkin method at $Ra = 700$ and $Ra = 10000$. The forcing terms $F_m$ are also plotted on the side.

## 5.5 Discussion

A closer look at the truncation schemes (figure 5.1.1b and 5.4.7b) provides some insight: in NG models, all forcing is embedded in the mean modes $(n = 0)$ via $F_m$ and thus any truncation confined to a finite box in the spectral space will only utilize the forcing contained in mean modes up to $M_{max}$. For $Ra = 700$ (figure 5.4.7a), once 50 percent of the modes in a $32 \times 48$ box are included, virtually no more forcing $(F_m)$ is left to be utilized. In contrast, at $Ra = 10000$ (figure 5.4.7b), all the forcing available in the $32 \times 48$ box is already employed at a 25 percent truncation and thus adding more modes will only add dissipation and the remainder of the forcing is left unused.

In FG models the modal distribution of forcing is entirely different. There, all forcing is left in the form of unstable modes, precisely those with the highest precedence in the standard truncation scheme. Hence the divergence of extreme truncations and the sharp fall in the Nusselt numbers produced by the smallest non-diverging truncations.

It is clear then that the finite-box constraint in our truncation scheme needs to be abandoned and the linear coefficient should be relied on as the sole determinant of a mode's precedence. This will guarantee, in both models, that an increasing sequence of truncations will exhaust all available forcing before the asymptotic trend toward the inclusion of high-frequency, high-dissipation modes sets in.

Figure 5.5.1 shows the results of another set of FG models, this time truncated in an "unconstrained" fashion. By that we mean that lying within a box of predefined size in the spectral space is no longer a required condition for the inclusion of a mode. The truncations are determined solely based on the linear coefficients as a criterion for modal precedence. Clearly, a more rapid and uniform convergence to the exact value is observed suggesting that

Figure 5.5.1: FG models at $Ra = 10^4$ with unconstrained truncations.

the anomalies in 5.4.5b and 5.4.6b are symptomatic of the "constrained" truncations.

It is only natural to surmise that this change to the truncation scheme should produce a similarly pronounced improvement in the performance of the NG models as well. However, at the time of writing, we do not have numerical evidence for this claim. Due to the high computational cost of the NG models, especially at high vertical wave numbers where the numerical eigenfunctions have to be resolved with rather high spatial resolutions, this task is best left to future investigations.

81

# CHAPTER VI

# Conclusions

The investigation the results of which we reported in this dissertation led to answers to a number of questions. First, our numerical simulations confirmed that in the turbulent regime of porous-medium convection, the coherent structures known as the minimal flow units are not only spatially repetitive units, but also the smallest "complete" dynamical units of the flow: we showed that a periodic box with the aspect ratio equal to the experimentally measured average width of the minimal flow unit is the smallest box capable of sustaining a flow with the "correct" vertical heat transport. Later, our minimal-flow-unit dynamical models were also able to reproduce the single mega-plume pair in that aspect ratio.

Secondly, we derived a class of Galerkin methods with a number of desirable properties, tailored to the problem of porous-medium convection. Our method is designed to model the boundary layers statically, relieving the dynamics of the task of resolving them. Furthermore, it isolates and separates all forcing in the form of constant inhomogeneities in a small localized set of the modal equations. This allows the linear truncation scheme to produce severely truncated working models by dispensing a balanced combination of forcing and dissipation. In contrast, Fourier-Galerkin models constructed using the same truncation scheme necessarily begin with all the forcing at once and thus diverge at severe truncations.

We are now also aware of the limitations of the method. Most importantly, the basis functions used need to be computed numerically in advance. This also means that the integrations and differentiations involved in the computation of the coefficients defining the ODEs have to be done numerically, adding tremendously to the computational cost of the method in comparison with the Fourier-Galerkin method.

In the beginning, the gap between the first two branches of the nonlinear spectrum and the higher branches (see figure 4.2.1) seemed to point to the unique dynamical role of modes chosen from those two branches. These modes indeed play a unique role in resolving fine structures near the boundary layers, but as we found later, unless many other modes from other branches are also included to resolve the bulk, we can not expect the models to satisfactorily reproduce the physical features of the flow. In particular, at high $Ra$, the forcing is spread over a large range of vertical wave numbers in the $n = 0$ modes. Thus, a converging model requires a large number of those modes.

Having understood the results presented in this report, we can propose a number of additional strategies for future exploration of the subject. The inhomogeneous forcing terms together with large self-inhibitions seem to force the large-$m$ modes from the $n = 0$ family into a more or less steady state. The energy spectra of the solutions (not reported here) paint a picture of those modes as steady sources distributing energy among other modes. Therefore, one may be able to model them statically, and incorporate their full forcing even in low truncations.

Another approach is to model the $n = 0$ modes separately from the rest, in the form of a one-dimensional PDE coupled to the ODEs, supplying them with forcing. This approach is currently being explored by our collaborators[6].

---

[6]Gregory Chini and Baole Wen at the University of New Hampshire.

It is our hope that this study paves the way for and inspires further research into reduced modeling of infinite-dimensional dynamical systems.

# APPENDICES

# APPENDIX A
# The inertial manifold

An elementary proof of existence for the global attractor requires the existence of a compact absorbing set for a dissipative evolution semigroup [3]. In some cases, one may additionally derive bounds on the Haussdorf dimension of the global attractor (and thus show that it is finite-dimensional) by following the evolution of arbitrary infinitesimal $n$-dimensional volumes under the evolution semigroup [2]. This is a considerably weaker result than the existence of a finite-dimensional exponentially absorbing smooth manifold containing the global attractor: the inertial manifold. It is the restriction to this manifold that renders the asymptotic dynamics finite-dimensional. Foias, Sell and Temam [8] first introduced this notion and provided a proof of existence applicable to a certain class of problems: for an evolution equation of the form

$$\frac{\mathrm{d}u}{\mathrm{d}t} + Au = F(u), \quad u \in H \tag{6.0.1}$$

where $H$ is a Hilbert space and $A : H \to H$ a positive linear operator, the so-called *strong squeezing property* guarantees the existence of the inertial manifold. This property is essentially the condition that two different initial states which are closer in their high-frequency component than the low-frequency component will remain so under the evolution, unless the two trajectories converge to one another exponentially. If $F : L^2 \to L^2$ is Lipschitz

continuous ($|F(u) - F(v)| \leq C_1 |u - v|$ for $u, v \in H$), it is then up to the linear operator $A$ to ensure that the magnitudes of the low and high components evolve in accordance with the squeezing property. This results in a condition relating the largest eigenvalue of $A$ in the low subspace $\lambda_N$, the smallest eigenvalue of $A$ in the high subspace $\lambda_{N+1}$, and the Lipschitz constant of $F$, $C_1$. Thus, a weak *spectral gap* criterion is obtained as a sufficient condition for the strong squeezing property:

$$\lambda_{N+1} - \lambda_N \geq 4C_1. \tag{6.0.2}$$

This is to say that there exists a large enough gap in the spectrum of $A$, and the smallest $N$ at which such a gap exists, determines the dimension of the minimal inertial manifold guaranteed to exist according to this method.

On the other hand, if we fail to demonstrate that $F : L^2 \to L^2$ is Lipschitz continuous, a more restrictive spectral gap criterion is required to produce the desired result. For instance, in the Navier-Stokes equations, the presence of the spatial derivatives in the nonlinear advection term prevents Lipschitz continuity in $L^2$. However, Lipschitz continuity may still be shown if $F$ is defined between so-called fractional power spaces, in which case the stronger spectral gap criterion is

$$\lambda_{N+1} - \lambda_N \geq 2C_1 \left( \lambda_{N+1}^\gamma + \lambda_N^\gamma \right) \tag{6.0.3}$$

where $\gamma$ is a positive number depending on the restricted domain and range of $F$. Now, it is required of the spectral gap to grow large, *soon enough* in order to counter the growth on the right hand side. This condition is not known to hold for $A = -\Delta$ (except in one spatial dimension where $\lambda_n \sim n^2$) and thus a proof of existence for the inertial manifold of the Navier-Stokes equations is currently not known.

# APPENDIX B
# Centro-symmetry

Graham and Steen [21] consider a $1 \times 1$ box with Neumann boundary conditions for temperature on the sidewalls and centro-symmetry (technically, antisymmetry). Because of the Neumann boundary condition, we can place such a field next to its reflection (across one of the side walls) and obtain a solution that is

1. periodic in $x$ over $[0, 2]$,

2. reflectionally symmetric about $x = 1$,

3. centro-symmetric within each of the two $1 \times 1$ sub-boxes,

4. differentiable everywhere including at $x = 1$,

and therefore consistent with the overall constraints and boundary conditions of our nonlinear model.

We now ask how we can enforce all these constraints on the solutions of our own nonlinear model. We expand the temperature field as follows:

$$\theta(x, z, t) = \sum_{m=0}^{\infty} \sum_{n=-\infty}^{+\infty} a_{mn}(t)\Theta_{mn}(z)e^{iknx} \tag{6.1.4}$$

for $x \in [0, 2]$ and $z \in [0, 1]$.

In order to enforce reflectional symmetry at $x = 1$, we demand that $\theta(x, z, t) = \theta(2 - x, z, t)$ for all $x, z, t$. Therefore,

$$\sum_{m,n} a_{mn}(t) \Theta_{mn}(z) \left[ e^{inkx} - e^{ink(2-x)} \right] = 0 \tag{6.1.5}$$

$$\sum_{m,n} 2i \sin(nkx) a_{mn}(t) \Theta_{mn}(z) = 2i \sum_{m} \sum_{n=1}^{\infty} \sin(nkx) \Theta_{mn}(z) \left[ a_{mn} - a_{mn}^* \right] \tag{6.1.6}$$

$$= -4 \sum_{m} \sum_{n=1}^{\infty} \sin(nkx) \Theta_{mn}(z) \operatorname{Im} a_{mn}(t) \tag{6.1.7}$$

$$= 0. \tag{6.1.8}$$

Therefore, $a_{mn}$ has to be real for all $m, n$.

The next constraint is centro-symmetry. To enforce this constraint, we introduce the following transformation:

$$x' = x - \frac{1}{2}, \quad z' = z - \frac{1}{2} \tag{6.1.9}$$

which translates the origin of the coordinate system to $(\frac{1}{2}, \frac{1}{2})$. Define

$$\tilde{\Theta}_{mn}(z') = \Theta_{mn}(z) = \Theta_{mn}(z' + \frac{1}{2}) \tag{6.1.10}$$

and write

$$e^{inkx} = e^{ink(x' + \frac{1}{2})} \tag{6.1.11}$$

$$= e^{in\frac{k}{2}} e^{inkx'}. \tag{6.1.12}$$

Thus, the temperature field written in terms of $x', z', t$ will take the form:

$$\theta(x', z', t) = \sum_m \sum_n a_{mn}(t)\tilde{\Theta}_{mn}(z')e^{in\frac{k}{2}}e^{inkx'} \tag{6.1.13}$$

$$= \sum_m a_{m0}\tilde{\Theta}_{m0}(z') + \sum_m \sum_{n=1}^{\infty} 2a_{mn}(t)\tilde{\Theta}_{mn}(z')\cos(nkx' + \frac{n}{2}k). \tag{6.1.14}$$

With this form of the expansion, we can now enforce centro-symmetry by demanding that

$$\sum_m a_{m0}\tilde{\Theta}_{m0}(z') + \sum_m \sum_{n=1}^{\infty} 2a_{mn}(t)\tilde{\Theta}_{mn}(z')\cos(nkx' + \frac{n}{2}k) \tag{6.1.15}$$

$$= -\sum_m a_{m0}\tilde{\Theta}_{m0}(-z') - \sum_m \sum_{n=1}^{\infty} 2a_{mn}(t)\tilde{\Theta}_{mn}(-z')\cos(-nkx' + \frac{n}{2}k). \tag{6.1.16}$$

As a result of the symmetry properties of $\Theta_{mn}(z')$, namely that

$$\Theta_{mn}(-z) = \begin{cases} \Theta_{mn}(z) & m \text{ odd} \\ \\ -\Theta_{mn}(z) & m \text{ even,} \end{cases} \tag{6.1.17}$$

this condition becomes:

$$0 \equiv \sum_{m\,\text{even}} 2a_{m0}(t)\tilde{\Theta}_{m0}(z') \tag{6.1.18}$$

$$+ 2\sum_{m\,\text{even}}\sum_{n=1}^{\infty} a_{mn}(t)\tilde{\Theta}_{mn}(z')\left[\cos(nkx' + \frac{n}{2}k) + \cos(-nkx' + \frac{n}{2}k)\right] \tag{6.1.19}$$

$$+ 2\sum_{m\,\text{odd}}\sum_{n=1}^{\infty} a_{mn}(t)\tilde{\Theta}_{mn}(z')\left[\cos(nkx' + \frac{n}{2}k) - \cos(-nkx' + \frac{n}{2}k)\right]. \tag{6.1.20}$$

Using the trigonometric identities

$$\cos(a+b) + \cos(-a+b) = 2\cos(a)\cos(b) \tag{6.1.21}$$

$$\cos(a+b) - \cos(-a+b) = -2\sin(a)\sin(b) \tag{6.1.22}$$

we have:

$$0 \equiv \sum_{m\,\text{even}} 2a_{m0}(t)\tilde{\Theta}_{m0}(z') \tag{6.1.23}$$

$$+ 4 \sum_{m\,\text{even}} \sum_{n=1}^{\infty} a_{mn}(t)\cos(\frac{n}{2}k)\tilde{\Theta}_{mn}(z')\cos(nkx') \tag{6.1.24}$$

$$- 4 \sum_{m\,\text{odd}} \sum_{n=1}^{\infty} a_{mn}(t)\sin(\frac{n}{2}k)\tilde{\Theta}_{mn}(z')\sin(nkx'). \tag{6.1.25}$$

Since $k = \pi$, we have

$$\cos(\frac{n}{2}k) = \begin{cases} 0 & n\,\text{odd} \\ \text{non-zero} & n\,\text{even} \end{cases} \tag{6.1.26}$$

and

$$\sin(\frac{n}{2}k) = \begin{cases} \text{non-zero} & n\,\text{odd} \\ 0 & n\,\text{even.} \end{cases} \tag{6.1.27}$$

Consequently, by the completeness of the bases $\tilde{\Theta}_{mn}(z')\cos(nkx')$ and $\tilde{\Theta}_{mn}(z')\sin(nkx')$, the only way the expression can be identically zero is if $a_{mn}(t) \equiv 0$ whenever the accompanying

trigonometric coefficient ( $\cos(\frac{n}{2}k)$ or $\sin(\frac{n}{2}k)$) is non-zero. In other words, we demand that

$$a_{mn}(t) = \begin{cases} 0 & m+n \ \text{even} \\ \text{Nonzero} & \text{otherwise} \end{cases} \tag{6.1.28}$$

An identical analysis leads to an identical result for the Fourier-Galerkin equations.

# APPENDIX C
# Temperature field snapshots from NG models



(a) $Ra = 700$, low resolution.

(b) $Ra = 700$, high resolution.

(c) $Ra = 1500$, low resolution.

(d) $Ra = 1500$, high resolution.

(e) $Ra = 10000$, low resolution.

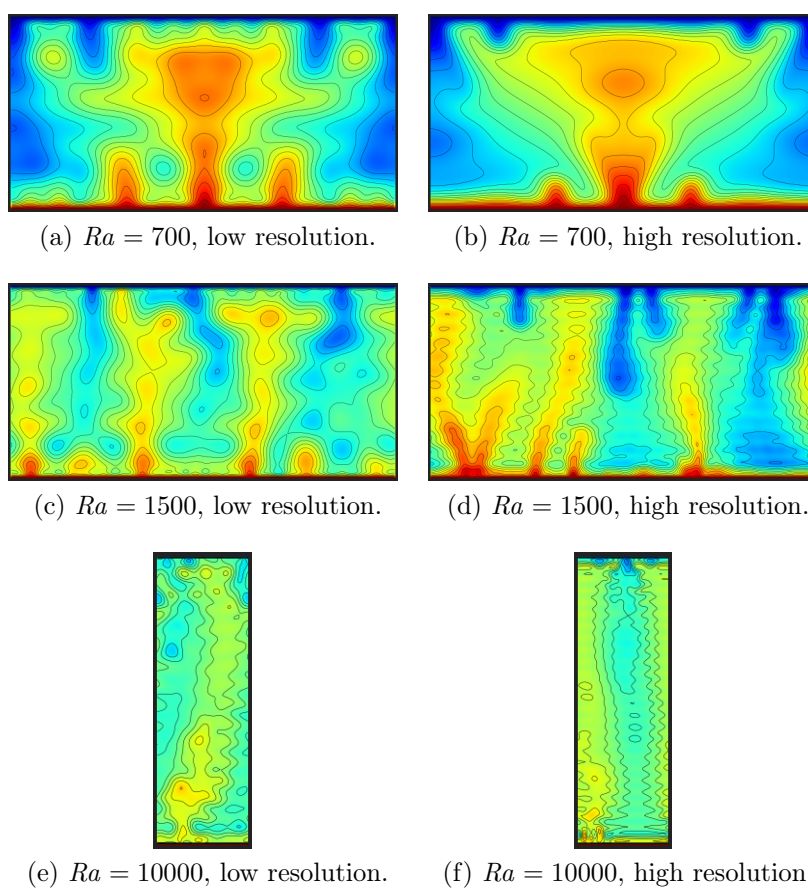(f) $Ra = 10000$, high resolution.

Figure 6.2.2: Examples of the temperature fields obtained from non-linear Galerkin models.

In this appendix, we present examples of the temperature field obtained from low and high-resolution non-linear Galerkin simulations. The low-resolution examples were chosen among

the lowest truncations capable of reproducing the main large-scale qualitative features of the flow. The "high resolution" examples were taken from the highest-resolution models simulated for each $Ra$. Not all are necessarily converged, but all demonstrate the large-scale qualitative features of the flow faithfully and robustly.

# Bibliography

[1] D. R. Hewitt, J. A. Neufeld, and J. R. Lister, "Ultimate Regime of High Rayleigh Number Convection in a Porous Medium," *Phys. Rev. Lett.*, vol. 108, p. 224503, May 2012.

[2] C. R. Doering, J. D. Gibbon, D. D. Holm, and B. Nicolaenko, "Exact Lyapunov dimension of the universal attractor for the complex Ginzburg-Landau equation," *Physical review letters*, vol. 59, no. 26, pp. 2911–2914, 1987.

[3] J. C. Robinson, *Infinite-dimensional dynamical systems: an introduction to dissipative parabolic PDEs and the theory of global attractors*, vol. 28. Cambridge University Press, 2001.

[4] R. Temam, *Infinite dimensonal dynamical systems in mechanics and physics*, vol. 68. Springer, 1997.

[5] C. Foias, O. Manley, and R. Temam, "Attractors for the Bénard problem: existence and physical bounds on their fractal dimension," *Nonlinear Analysis*, vol. 11, no. 8, pp. 939–967, 1987.

[6] M. A. Efendiev, J. Fuhrmann, and S. V. Zelik, "The long-time behaviour of the thermoconvective flow in a porous medium," *Mathematical methods in the applied sciences*, vol. 27, no. 8, pp. 907–930, 2004.

[7] M. A. Efendiev and S. V. Zelik, "Finite-and infinite-dimensional attractors for porous media equations," *Proceedings of the London Mathematical Society*, vol. 96, no. 1, pp. 51–77, 2008.

[8] C. Foias, G. Sell, and R. Temam, "Inertial manifolds for nonlinear evolutionary equations," *Journal of Differential Equations*, vol. 73, no. 2, pp. 309–353, 1988.

[9] J. P. Boyd, *Chebyshev and Fourier spectral methods*. Dover publications, 2001.

[10] C. Canuto, M. Y. Hussaini, A. Quarteroni, and T. A. Zang, *Spectral methods in fluid dynamics*. Springer-Verlag, 1988.

[11] P. Holmes, J. Lumley, and G. Berkooz. Turbulence, coherent structures, dynamical systems and symmetry, Cambridge University Press, 1998.

[12] T. R. Smith, J. Moehlis, and P. Holmes, "Low-dimensional modelling of turbulence using the proper orthogonal decomposition: a tutorial," *Nonlinear Dynamics*, vol. 41, no. 1, pp. 275–307, 2005.

[13] J. Bailon-Cuba and J. Schumacher, "Low-dimensional model of turbulent Rayleigh-Bénard convection in a Cartesian cell with square domain," *arXiv preprint arXiv:1106.4157*, 2011.

[14] D. A. Nield and A. Bejan, *Convection in porous media*. Springer, 2006.

[15] B. Metz, O. Davidson, H. De Coninck, M. Loos, and L. Meyer, "IPCC special report on carbon dioxide capture and storage," tech. rep., 2005.

[16] T. Dubois, F. Jauberteau, and R. Temam, *Dynamic multilevel methods and the numerical simulation of turbulence*. Cambridge University Press, 1999.

[17] J. W. Elder, "Steady free convection in a porous medium heated from below," *J. Fluid Mech*, vol. 27, no. 1, pp. 29–48, 1967.

[18] M. D. Shattuck, R. P. Behringer, G. A. Johnson, and J. G. Georgiadis, "Convection and flow in porous media. Part 1. Visualization by magnetic resonance imaging," *Journal of Fluid Mechanics*, vol. 332, pp. 215–246, 1997.

[19] J. Otero, L. A. Dontcheva, H. Johnston, R. A. Worthing, A. Kurganov, G. Petrova, and C. R. Doering, "High-Rayleigh-number convection in a fluid-saturated porous layer," *Journal of Fluid Mechanics*, vol. 500, no. 1, pp. 263–281, 2004.

[20] M. D. Graham and P. H. Steen, "Strongly interacting travelling waves and quasiperiodic dynamics in porous medium convection," *Physica D: Nonlinear Phenomena*, vol. 54, no. 4, pp. 331–350, 1992.

[21] M. D. Graham and P. H. Steen, "Plume formation and resonant bifurcations in porous-media convection," *Journal of Fluid Mechanics*, vol. 272, pp. 67–90, 1994.

[22] P. Constantin and C. R. Doering, "Heat transfer in convective turbulence," *Nonlinearity*, vol. 9, no. 4, pp. 1049–1060, 1996.

[23] C. R. Doering and P. Constantin, "Variational bounds on energy dissipation in incompressible flows: Shear flow," *Physical Review E*, vol. 49, no. 5, p. 4087, 1994.

[24] C. R. Doering and P. Constantin, "Bounds for heat transport in a porous layer," *Journal of Fluid Mechanics*, vol. 376, no. 1, pp. 263–296, 1998.

[25] E. R. Lapwood, "Convection of a fluid in a porous medium," *Proceedings of the Cambridge*, 1948.

[26] S. Backhaus, K. Turitsyn, and R. E. Ecke, "Convective Instability and Mass Transport of Diffusion Layers in a Hele-Shaw Geometry," *Phys. Rev. Lett.*, vol. 106, p. 104501, Mar 2011.

[27] G. B. Arfken and H. J. Weber, *Mathematical methods for physicists*. Academic Press, 2001.

[28] W. V. R. Malkus, "The heat transport and spectrum of thermal turbulence," *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*, vol. 225, no. 1161, pp. 196–212, 1954.

[29] L. N. Howard, "Heat transport by turbulent convection," *Journal of Fluid Mechanics*, vol. 17, no. 03, pp. 405–432, 1963.

[30] L. N. Howard, "Convection at high Rayleigh number (Thermal convection in horizontally infinite layer of fluid confined between rigid heat conducting plates driven by temperature difference)," *Applied Mechanics, Proceedings of the Eleventh International Congress of Applied Mechanics, Munich, West Germany; 30 Aug.-5 Sept.1964*, pp. 1109–1115, 1966.

[31] R. N. Horne and M. J. O'sullivan, "Origin of oscillatory convection in a porous medium heated from below," *Phys. Fluids;(United States)*, vol. 21, no. 8, 1978.

[32] P. Castiglione, M. Falcioni, A. Lesne, and A. Vulpiani, *Chaos and coarse graining in statistical mechanics*. Cambridge University Press, 2008.

[33] L. N. Trefethen, *Spectral methods in MATLAB*, vol. 10. Society for Industrial Mathematics, 2000. This an example note.

[34] R. LeVeque, *Finite Difference Methods for Ordinary and Partial Differential Equations: Steady-State and Time-Dependent Problems (Classics in Applied Mathemat)*. Philadelphia, PA, USA: Society for Industrial and Applied Mathematics, 2007.

[35] E. A. Coutsias, T. Hagstrom, and D. Torres, "An Efficient Spectral Method for Ordinary Differential Equations with Rational Function Coefficients," *Mathematics of Computation*, vol. 65, no. 214, pp. –611, 1996.

[36] T. R. Smith, J. Moehlis, and P. Holmes, "Low-dimensional models for turbulent plane Couette flow in a minimal flow unit," *Journal of Fluid Mechanics*, vol. 538, pp. 71–110, 2005.

[37] G. E. Andrews, R. Askey, and R. Roy, *Special Functions, volume 71 of Encyclopedia of Mathematics and its Applications.* Cambridge University Press Cambridge, 1999.

[38] C. R. Doering and P. Constantin, "Energy dissipation in shear driven turbulence," *Physical review letters*, vol. 69, no. 11, pp. 1648–1651, 1992.

[39] G. Chini, N. Dianati, Z. Zhang, and C. Doering, "Low-dimensional models from upper bound theory," *Physica D: Nonlinear Phenomena*, vol. 240, no. 2, pp. 241–248, 2011.

[40] B. Wen, N. Dianati, E. Lunasin, G. P. Chini, and C. R. Doering, "New upper bounds and reduced dynamical modeling for Rayleigh-Bénard convection in a fluid saturated porous layer," *Communications in Nonlinear Science and Numerical Simulation*, vol. 17, no. 5, pp. 2191–2199, 2012.

[41] G. P. Chini, B. Wen, N. Dianati, and C. R. Doering, "Computational approaches to aspect-ratio-dependent upper bounds and heat flux in porous medium convection," *submitted*.

[42] M. E. J. Newman, *Computational Physics.* CreateSpace Independent Publishing Platform, 2012.