

Deep learning for surrogate modeling of two-dimensional mantle convection

Siddhant Agarwal,^{1,2,*} Nicola Tosi,¹ Pan Kessel,^{2,†} Doris Breuer,¹ and Grégoire Montavon^{2,†}

¹*Planetary Physics, Institute of Planetary Research,
German Aerospace Center (DLR), Berlin, Germany*

²*Machine Learning Group, Berlin Institute of Technology, Berlin, Germany*

(Dated: November 8, 2021)

arXiv:2108.10105v2 [astro-ph.EP] 5 Nov 2021

Abstract

Mantle convection, the buoyancy-driven creeping flow of silicate rocks in the interior of terrestrial planets like Earth, Mars, Mercury and Venus, plays a fundamental role in the long-term thermal evolution of these bodies. Yet, key parameters and initial conditions of the partial differential equations governing mantle convection are poorly constrained. This often requires a large sampling of the parameter space to determine which combinations can satisfy certain observational constraints. Traditionally, 1D models based on scaling laws used to parameterized convective heat transfer, have been used to tackle the computational bottleneck of high-fidelity forward runs in 2D or 3D. However, these are limited in the amount of physics they can model (e.g. depth dependent material properties) and predict only mean quantities such as the mean mantle temperature. A recent machine learning study has shown that feedforward neural networks (FNN) trained using a large number of 2D simulations can overcome this limitation and reliably predict the evolution of entire 1D laterally-averaged temperature profile in time for complex models. We now extend that approach to predict the full 2D temperature field, which contains more information in the form of convection structures such as hot plumes and cold downwellings. Using a dataset of 10,525 two-dimensional simulations of the thermal evolution of the mantle of a Mars-like planet, we show that deep learning techniques can produce reliable parameterized surrogates (i.e. surrogates that predict state variables such as temperature based only on parameters) of the underlying partial differential equations. We first use convolutional autoencoders to compress the size of each temperature field by a factor of 142 and then use FNN and long-short term memory networks (LSTM) to predict the compressed fields. On average, the FNN predictions are 99.30% and the LSTM predictions are 99.22% accurate with respect to unseen simulations. Proper orthogonal decomposition (POD) of the LSTM and FNN predictions shows that despite a lower mean relative accuracy, LSTMs capture the flow dynamics better than FNNs. When summed, the POD coefficients from FNN predictions and from LSTM predictions amount to 96.51% and 97.66% relative to the coefficients of the original simulations, respectively.

I. INTRODUCTION

Studying the long-term thermal evolution of terrestrial planets like Earth, Venus, Mercury and Mars requires detailed modelling of thermal convection in their rocky mantles (e.g., [1]). Similar

* Funded by HEIBRiDS Graduate School for Data Science; agsiddhant@gmail.com

† Funded by BIFOLD

to the flow of crystalline ice in glaciers, mantle convection is a form of sub-solidus convection. Mantle rocks at high temperature and pressure, but still well below their melting point, behave like a highly viscous fluid over geological time scales (millions to billions of years), largely in response to thermal and compositional buoyancy forces. Mantle convection is typically modelled using fluid dynamics codes (e.g., [2–5]) that numerically solve the non-linear partial differential equations (PDEs) of mass, momentum and energy conservation governing the creeping flow (i.e. with negligible inertia) of silicate rocks subject to basal heating from the metallic core and internal heating due to the decay of radioactive elements. Unfortunately, key model parameters (such as the rock viscosity, or the amount of internal heating) and initial conditions (such as the initial mantle and core temperatures), that are inputs to these forward numerical models, are poorly constrained. Thus, one typically chooses and tests a large number of different parameter values (within a reasonable range) to observe how these affect the outputs of the simulations. The outputs can be processed to arrive at various quantities of interest such as surface heat flux, amount of thermal contraction or expansion, duration and timing of volcanism, or volume of produced crustal material. To some extent, these quantities are “observables” that can be inferred from geophysical and geochemical data delivered by planetary space missions. In turn, they provide fundamental constraints on the convective evolution of terrestrial bodies (see [6] for a recent review about this topic).

This approach of testing several different values of parameters, however, suffers from a bottleneck imposed by the computational cost of the 2D and 3D forward models. No matter whether one approaches this issue from the perspective of an inverse problem (inferring the model parameters given the observables) or of a forward problem (calculating the observables given the model parameters), it is often impractical to run several thousands of simulations to determine which parameters and combinations thereof can satisfy a set of given observational constraints.

A number of inverse-problem studies have attempted to overcome this computational bottleneck of expensive simulations, ranging from using modified Markov Chain Monte Carlo (MCMC) methods (e.g. [7, 8]), all the way to completely bypassing MCMC methods and directly learning the mapping between parameters and observables from simulations run prior to the inversion using Mixture Density Networks (MDN) (e.g., [9–13]).

Machine learning (ML) methods such as MDN can also be used to learn highly non-linear forward mappings from parameters to observables. These can preserve some physical insights into the flow being modelled (such as convection patterns) in contrast to purely statistical inferences

made under an inverse formulation. Traditionally, several mantle convection studies have employed “scaling laws”, which parameterize the heat flow out of the mantle (quantified by the Nusselt number (Nu), or the non-dimensional heat flux), in terms of the vigor of convection (quantified by the Rayleigh number (Ra)) (e.g., [14–17]). These scaling laws are then used in the frame of one-dimensional, spherically-symmetric models to advance the mean mantle temperature (and a number of associated quantities) in time by solving two ordinary differential equations governing the global energy balance of the mantle and core (e.g., [18–25]).

On the one hand, the use of scaling laws makes models of planetary evolution computationally very efficient as these only require the solution of ordinary differential equations. On the other hand, scaling laws are limited in the amount of physics they can capture. For example, solid-solid phase-transitions or the pressure-dependence of the viscosity and other thermal and transport properties such as thermal expansivity and conductivity can hardly be taken into account in scaling laws for heat transfer. Additionally, they allow one only to predict the evolution of global quantities such as the surface heat flux or the mean mantle temperature, but do not provide any insight into the spatial and temporal variability of mantle flow. Recently, [26] showed that some of these limitations can be overcome by using machine learning. They showed that feedforward neural networks (FNN) can be used for predicting the surface heat flux and mean temperature of steady-state simulations, given parameters such as Ra , the core-to-planet radius ratio, and mode of convection, i.e. mobile lid, characterizing Earth’s plate tectonics, or stagnant lid, characterizing the other terrestrial bodies of the solar system (see also Sec. II). [27] demonstrated that limitations related to taking into account complex physics and the time-variability of the heat transfer can be overcome through FNNs properly trained with a large set of 2D dynamic simulations. By using multiple parameters such as e.g. mantle reference viscosity (related to Ra) and activation volume and activation energy of diffusion creep rheology (controlling the temperature and pressure dependence of the viscosity) as inputs to the FNN, it is possible to directly predict the thermal evolution of the entire horizontally-averaged 1D temperature profile of the mantle in time with a mean accuracy of 99.7%. Another interesting study uses Generative Adversarial Networks to reconstruct missing plate boundaries derived from horizontal divergence maps of a steady-state 3D convection model [28].

In this paper, we build upon [27]. While the full 1D temperature profile already provides a lot more information than simply the surface heat flux and mean temperature and is also beyond what can be expected to be retrieved for planets like Mars through future planetary missions, it still lacks the rich convection structures such as plumes and downwellings that are delivered by 2D or 3D

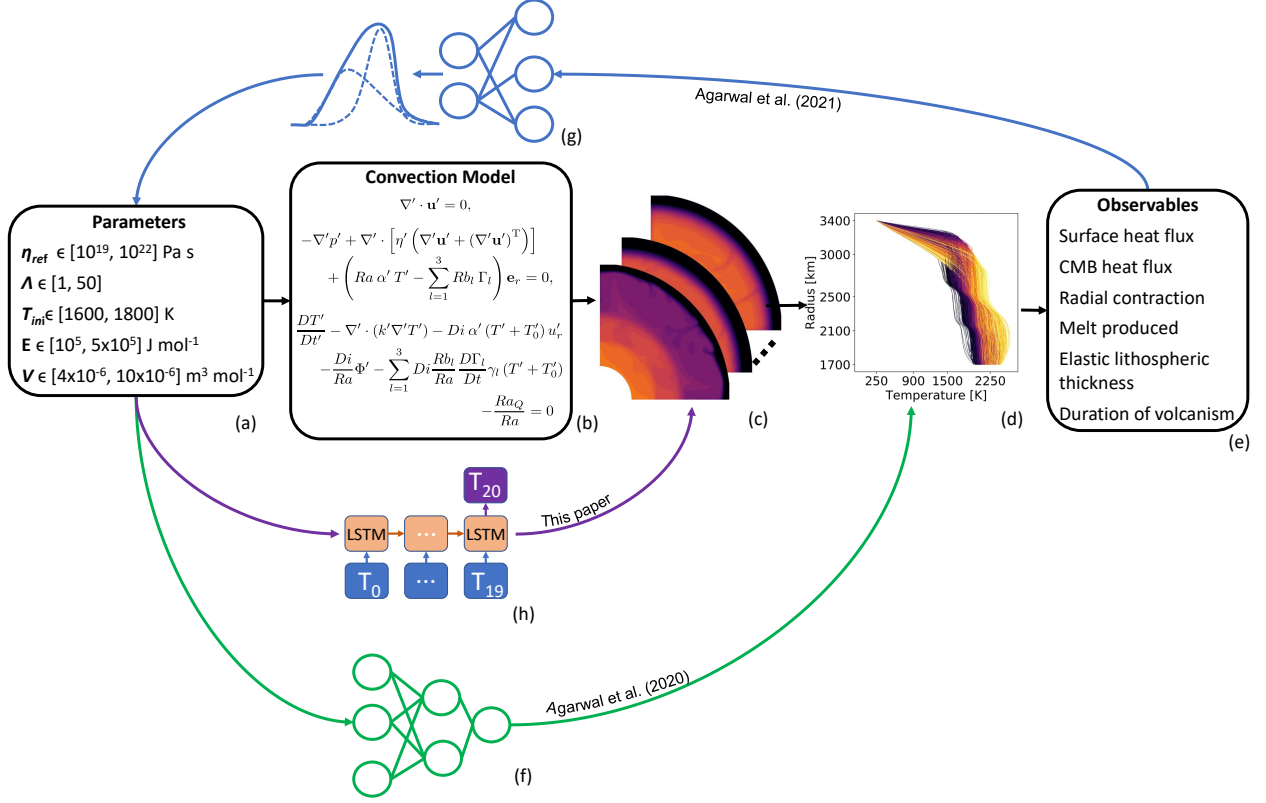


FIG. 1. The context for this study. (a) Typically, a mantle convection study starts by randomly drawing input parameters from a flat distribution and then feeding them to the forward models. (b) The PDEs are solved using dedicated mantle convection codes. (c) The outputs of the simulations can be processed to arrive at certain lower-dimensional observables such as (d) the horizontally-averaged 1D temperature profiles or (e) more global quantities such as the surface heat flux, radial contraction, duration of volcanism, etc. ML methods have been shown to work well for these low-dimensional observables, both - (f) in a forward study ([27]) and (g) an inverse study ([13]). In this work, we demonstrate that (h) a surrogate can model 2D mantle convection using deep learning.

simulations. In fact, [13] demonstrated that using the same setup for a Mars-like planet as in [27], the surface heat flux provides only a loose constraint for parameters governing mantle convection and even knowledge of the entire 1D temperature profile is insufficient for constraining certain parameters (e.g. the activation volume of diffusion creep rheology, which controls the pressure dependence of the viscosity). In other words, more information (such as horizontal variations of the temperature field) might hold the key to placing tighter constraints on certain thermal-evolution parameters. Hence, it is desirable to construct surrogate models capable of predicting

more information about a convecting mantle.

Here, we show how deep learning (DL) can be leveraged to directly predict surrogates of the 2D temperature fields from five key-parameters. Figure 1 shows the concept and general context within which this study fits. We sample the model parameters from a flat distribution with a broad, yet reasonable range and input these to the convection simulations. The five key parameters are: (1) reference mantle viscosity (η_{ref}), which is the viscosity attained at a given reference depth and reference temperature; (2) activation energy of the diffusion creep (E), which controls the temperature dependence of the viscosity, and (3) activation volume of diffusion creep (V) controlling the degree to which viscosity depends on pressure; (4) enrichment factor (Λ), which determines the proportion of the heat-producing radiogenic elements extracted from the convecting mantle upon melting and enriched in the crust; (5) the initial temperature of the mantle (T_{ini}). These five parameters typically have the largest influence on the forward model for the thermal evolution of terrestrial planets that we use (e.g. [29], [30]).

To our knowledge, this is the first time a parameterized 2D surrogate is proposed in the context of mantle convection (see [31] for a recent review of applications of data science methods in geodynamics). We use the term parameterized to stress that we are interested in predicting a variety of mantle flows based on different combinations of input parameters and not, for example, given a certain amount of time-steps of a single simulation with fixed parameters, predicting the subsequent time-steps. [32] demonstrated that direct numerical simulation (DNS) of two-dimensional turbulent Rayleigh-Bénard convection can be modelled using reservoir computing. In the above study, the time-steps of the same, single simulation are split into training and test sets. In contrast, we are interested in modelling all the time-steps of different simulations, which is made possible by the relatively low computational cost of each simulation (20 to 500 CPU hours depending on the combination of parameters).

In general, machine learning for predicting flows is an active research area. Particularly worth highlighting is the seminal paper by [33], which showed that the PDEs can be embedded in the loss function through automatic differentiation of state variables with respect to spatial and temporal coordinates. However, here we stick to a purely data-driven approach, reserving the application of methods such as those in [33] for future research. A notable example of such a purely data-driven approach was shown to be effective in capturing the dynamics of 3D turbulence by [34]. They showed that the velocity fields can be compressed using convolutional autoencoders [35]. Convolutional autoencoders successively down-size the original field (or image) into a bottleneck,

from where they are reconstructed back to the original size. In this way, the dimensionality of the original high-resolution fields can be decreased and made more computationally efficient to work with. [34] then predicted these compressed time-steps using a convolutional long short-term memory (LSTM) network [36], which by providing a mechanism to relate the time-steps of a simulation, allows one to learn the attractor for the underlying dynamics. Just as [32], [34] split the time-steps of the same simulation into training and test sets. We show that this approach can be adapted to our purposes, so that we can predict all the time-steps for any simulation in our data manifold, given just a set of five parameters. Somewhat closer to our purposes, [37] used a convolutional encoder-decoder architecture for predicting pressure and temperature fields around airfoils, given the spatial grid as well as two additional parameters (angle of attack and Reynolds number). We refer to [38] for an overview of several ML techniques that have been used for prediction, dimensionality reduction and optimization and control in fluid dynamics.

The outline of the paper is as follows. In the Sec. II, we briefly outline the setup of the simulations of a Mars-like planet. Then, in Sec. III, we present how we compress the temperature fields obtained from the mantle convection simulations. We follow up with Sec. IV, where we predict the compressed temperature fields using two different ML algorithms - FNN and LSTM. In the same section, we further delve into the differences in the predictions of the two algorithms by analyzing them from the lens of Proper Orthogonal Decomposition (POD) or its ML equivalent, Principal Component Analysis. We then conclude the paper by offering some potentially interesting follow-ups to this work.

II. DATASET OF MANTLE CONVECTION SIMULATIONS

We employ a dataset consisting of 10,525 simulations of the thermal evolution of a Mars-like planet run on a 2D quarter-cylindrical grid (Fig. 2). A detailed description of the setup is provided in Appendix A. Here, we summarize the main features of the model.

We model the mantle as a viscous fluid with negligible inertia (i.e. with infinite Prandtl number or, equivalently, undergoing Stokes flow). We consider a pressure- and temperature-dependent Newtonian rheology [39], which is calculated using the Arrhenius law for diffusion creep, whose dimensional form reads

$$\eta(T, P) = \eta_{\text{ref}} \exp\left(\frac{E + PV}{T} - \frac{E + P_{\text{ref}}V}{T_{\text{ref}}}\right). \quad (1)$$

The reference viscosity η_{ref} is attained at reference temperature $T_{\text{ref}} = 1600$ K and reference pressure $P_{\text{ref}} = 3$ GPa, respectively. P is the hydrostatic pressure, E is the activation energy, and V is the activation volume.

Mars, in contrast to the Earth, but like Mercury, the Moon and, at least at present-day, Venus, operates in a so-called stagnant-lid convection mode (see e.g. [6]). The strong temperature dependence of the viscosity causes the relatively cold upper part of the mantle to develop a high-viscosity layer – the stagnant lid. Such a stiff layer remains immobile during the entire evolution of the planet (although its thickness can decrease or increase in response to heating or cooling of the mantle). The stagnant lid insulates the interior causing thermal convection to take place only in the mantle beneath it and to be largely driven by cold downwellings developing at its base (see Fig. 2). This mode of convection is remarkably different from the plate tectonic (or mobile-lid) mode of convection that characterizes the Earth. Cold tectonic plates are in fact an active part of the Earth’s convecting engine: by sinking into the mantle at subduction zones, they provide a strong cooling for the mantle and core with fundamental consequences for large-scale transport of materials in the deep interior and for the generation of the Earth’s magnetic field.

Within the relatively low pressure range of Mars mantle (~ 20 GPa), the degree of compressibility of mantle rocks is limited; the dissipation number (eq. A11) is quite small for Mars (~ 0.13 , significantly smaller than the Earth’s, which amounts to ~ 0.5). In this situation, it is appropriate to employ the so-called extended Boussinesq approximation (e.g., [40]). Like the standard Boussinesq approximation, the extended one assumes constant density everywhere except in the buoyancy term of the momentum equation (eq. A2). However, it further accounts for the effects of adiabatic compression/decompression and viscous dissipation in the conservation equation for the thermal energy (eq. A3).

The mantle temperature evolves in time due to the decay of radiogenic elements present in the mantle and due to the cooling of the core, which also provides a source of basal heat. Following a standard approach adopted in the mantle convection community, the core is simply considered as a homogeneous sphere of given density and heat capacity, whose mean temperature evolves at a rate imposed by the cooling rate of the mantle (see Fig. 2, eq. A21 and e.g., [18]).

Geological evidence suggests that the bulk of the crust of Mars formed very early in the evolution of the planet [41]. Similar to other studies (e.g., [30]), we thus assume that a crust of a fixed thickness has been present since the beginning of the evolution. This assumption implies that, from the beginning of each simulations, the mantle is partly depleted of radiogenic elements

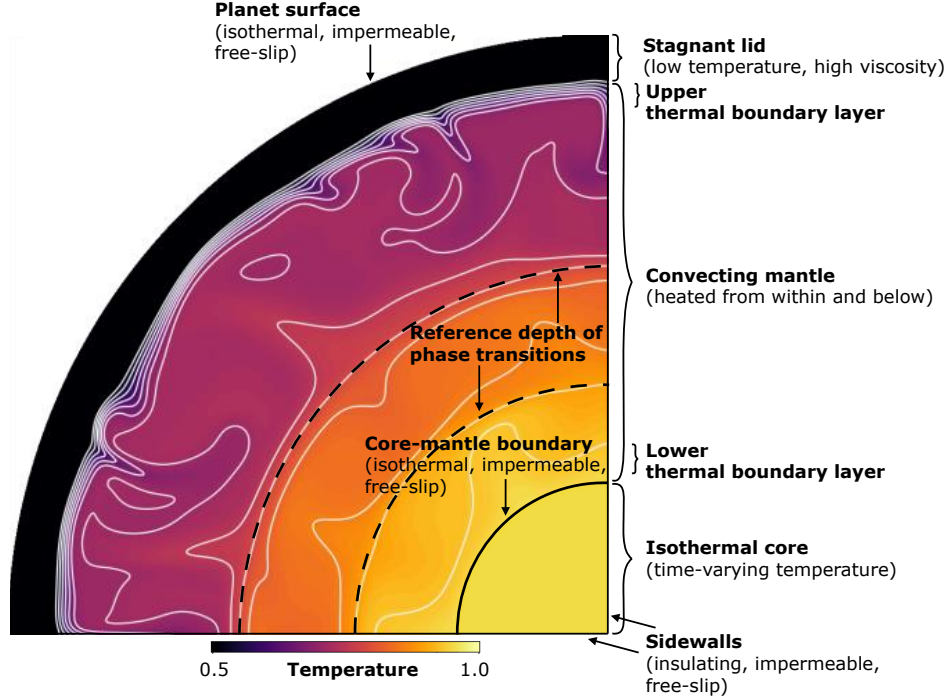


FIG. 2. Two-dimensional, quarter-cylindrical domain illustrating the main features of the employed mantle convection model. The mantle is colored according to the non-dimensional temperature field of one specific simulation of the dataset. The diagram depicts the smaller rescaled core (see Sec. A 8). Since a large part of the temperature variations occur across the stagnant-lid, the colorscale is truncated at 0.5 for ease of visualization. See text for more details and Appendix A for a complete model description.

(resulting in new new depleted mantle composition C_{depleted}) with respect to their primordial concentration (C_0), which we set according to the model of [42] as follows:

$$C_{\text{depleted}} = \frac{M_m C_0}{M_{\text{cr}} (\Lambda - 1) + M_m}, \quad (2)$$

where M_m and M_{cr} are the mass of the mantle and crust, respectively and Λ is the crustal enrichment factor. The rationale of this assumption is that upon partial melting of the mantle – a common event during the evolution of terrestrial bodies – radiogenic elements behave as incompatible elements, i.e. they tend to be enriched in the melt phase with respect to the solid phase. As a consequence, the production of crust, which results from mantle melting, melt migration toward the surface and solidification, causes a net depletion of radiogenic elements in the mantle. In practice, we reduce the primordial bulk abundances of radiogenic elements according to a crustal enrichment factor Λ . The mantle is further depleted in radiogenic heat-producing elements during the evolution

whenever local partial melting takes place according to the model described by [43] (see Sec. A 6).

We also consider the influence on the mantle flow of two major solid-solid phase-transitions using the standard phase-function approach of [44] adopted in mantle convection (Sec. A 5). We further assume that the coefficient of thermal expansion and the thermal conductivity depend on pressure and temperature as appropriate for silicate materials (see Sec. A 4 and [45]).

The simulations are initialized with a constant mantle temperature T_{ini} combined with upper and lower 300-km-thick thermal boundary layers. For all simulations, a small random perturbation is added to the temperature field to initiate convection. The initial mantle temperature has a strong effect on the early evolution of the planet. However, this initial condition becomes less important with time (after ~ 2 billion years) due to the “thermostat effect” [46] - the temperature dependence of the viscosity regulates the temperature of the mantle. On the one hand, when the mantle is hot, the viscosity decreases, leading to more vigorous convection and thereby, efficient cooling. On the other hand, when the mantle is cooler at latter stages in the evolution, it is more viscous and cools less efficiently. While in this dataset of forward models we only vary the initial mantle temperature, other parameters related to the initial conditions that can have a potential impact on the thermal evolution of the interior could also be considered. For example, the starting temperature at the core-mantle boundary has implications for the dynamics of the lower mantle and affects melt and magnetic field generation (e.g. [47]). The initial thickness of the thermal boundary layers, although irrelevant for the long-term evolution, influences the heat fluxes across the surface and CMB during the first few hundred million years. Additionally, an initial compositional stratification resulting from the crystallization of a primordial magma ocean, neglected in our isochemical models, could also have a significant influence on the evolution of the mantle and core ([48–50]).

As for the boundary conditions, all domain boundaries are impermeable and free-slip. The surface temperature is kept fixed at 250 K throughout the evolution. Latitudinal variations of the surface temperature, as on Mars ([51]) do not have a significant impact on the long-term evolution and large-scale dynamics of the planet ([52]) and thus can be safely neglected. The temperature of core-mantle boundary evolves as the core cools (eq. A21). There is no heat flux across the side walls of the computational domain, i.e. they are assumed to be insulating (see Fig. 2). Table II and Table III list all the fixed dimensional parameters that are shared by all simulations.

We ran several single-core simulations using the finite-volume code GAIA [5] with a grid resolution of 300 radial layers and 392 cells per layer. Even though 3251 out of 10,525 simulations did not reach the end time of 4.5 Gyr (i.e. time since formation of the planet until today), we use

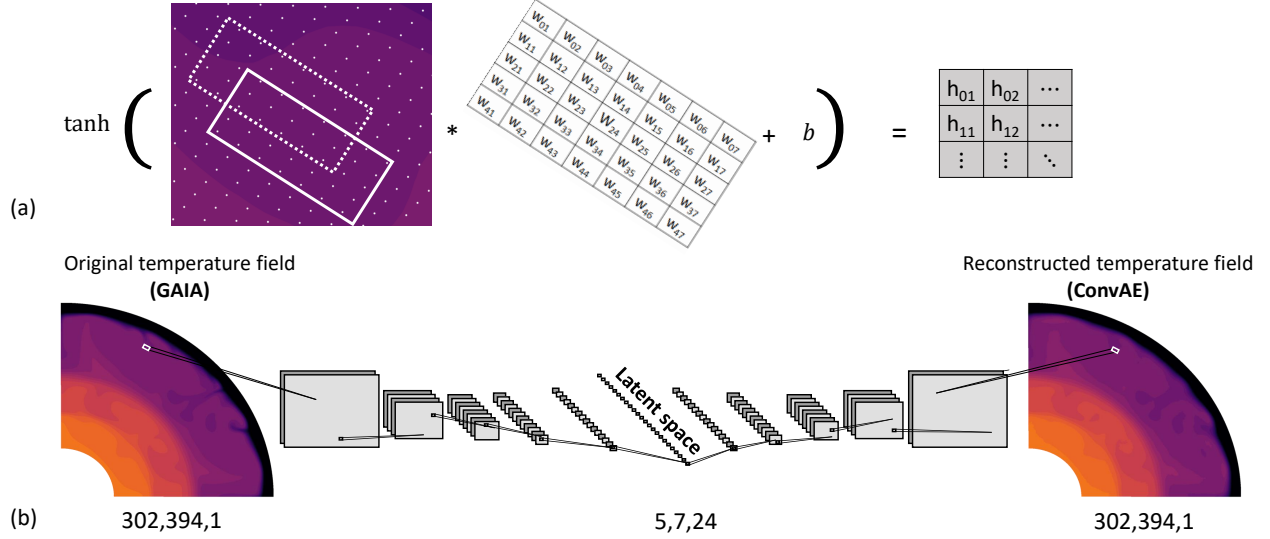


FIG. 3. Dimensionality reduction of the 2D temperature fields using convolutional autoencoders. (a) Filters with trainable parameters (\mathbf{w}) move across the computational domain with specified strides. After the convolution operation ($*$), the bias b is added to it before applying the activation function, resulting in the entries for the next hidden layer \mathbf{h} . (b) Several filters can be used to successively reduce the size of the original field ($302 \times 394 \times 1$) till a desired encoding or latent space representation is obtained ($5 \times 7 \times 24$). The convolutional autoencoder then reconstructs the compressed field back to its original dimensions ($302 \times 394 \times 1$) using deconvolution operations and optimizes the weights using the difference between the prediction and the ground truth.

all the time-steps available from all the simulations. Some combinations of parameters can make the system of PDEs too stiff to efficiently solve, as the time-stepping becomes increasingly small. This means that some simulations can take more than 50 times the average run time. In total, the dataset amounts to 10 TB and took approximately 2 million CPU hours. The challenging task of predicting the 2D temperature field for a variety of parameters necessitated this computational effort. The number of simulations used in this study is therefore almost a 100 times that of what a typical parameter study in mantle convection would use. The distribution of the parameters for all the simulations in the training, cross-validation and the test set is plotted in Fig. 16.

III. COMPRESSION OF TEMPERATURE FIELDS

The temperature fields alone account for approximately 1 TB. While most everyday computers cannot hold that much data in memory, one could still overcome the issue using a “data-generator”. A data-generator feeds the ML algorithm during training and/or inference by reading and storing only a limited number of examples from the disc at any given time. Nevertheless, these require careful programming to ensure that the CPU or GPU are constantly fed batches of data without having to wait for the next one. Furthermore, performing ML calculations on the full-sized fields would be slower than a compressed version of it and introduce a lot more trainable parameters which could increase the risk of over-parameterization. Hence, we decided to compress the temperature fields first and then scan for different ML algorithms and architectures that can help us predict this latent space representation.

Traditionally, linear reduced-order modelling techniques like proper orthogonal decomposition (POD) have been used for truncating high-fidelity simulations using the most dominant modes [53]. However, the orthonormal bases (typically obtained through singular value decomposition) of one simulation do not generalize well onto those of another and often require non-trivial basis interpolation (e.g., [54]). Recently, [34] demonstrated that convolutional autoencoders (ConvAE) [35] provide a powerful non-linear tool for compressing flow fields, bypassing the need for calculating POD modes.

Instead of fully connected dense layers like conventional autoencoders, ConvAE uses convolutional filters. As shown in Fig. 3(a), a filter with trainable weights \mathbf{w} moves across the state variable field (temperature) as specified by a hyperparameter called “stride”. A stride of 2, for example, means that the filters move two units (two numerical grid cells) horizontally and then when a row is completed, two units vertically. Fig. 3(a) shows a filter with height 5 and length 7 (also hyperparameters) convolving with the temperature field at strides of 2 in both x - and y -direction. We use $\tanh()$ as activation function, which, when applied to the sum of the bias and the convolution product, introduces non-linearity and returns the output for the next hidden layer which can then be convolved on. This process continues until the desired latent space representation is reached (Fig. 3(b)). Then, the compressed state can be successively restored to the original size with the help of another sequence of convolution filters, called deconvolutions in this context. Once the forward graph has been setup, one can then minimize the difference between the original and reconstructed 2D temperature field by back-propagating the derivative of the error with respect to

the network weights. [35] point out that a convolutional architecture offers two main benefits over the fully connected dense layers: (1) the 2D structures of the convolutional filters retain spatial correlations which would otherwise be lost in 1D dense layers and (2) ConvAE have significantly fewer trainable parameters due to the shared weights.

We use Keras, an API built on top of Tensorflow [55] for training our ConvAE. We split the simulations into 98%, 1% and 1% for training, validating and testing, respectively. Given the size of the entire training set, we only feed the GPU mini-batches of 16 temperature fields (i.e. time-steps of any simulation) during training. We use L2 regularization and early-stopping by manually monitoring the validation loss (mean-squared error) and the optimization is carried out using Adam [56].

Since, the computer cannot hold the entire training-set in memory, we use a data-generator. Using multi-processing built into Keras’ “fit-generator” for creating multiple batches in parallel along with multi-threading for populating each batch with the help of Joblib [57], it takes around 6 hours for one epoch to complete. After 5–10 epochs, acceptable results are obtained. Fig. 4(a) and 4(b) show results of reconstruction for two different examples in the test set using three different ConvAE architectures, which differ in the dimensions of the latent space. As it can be expected, the more the temperature fields are compressed, the less accurate the reconstruction is. In all the plots of the temperature field in this paper, we clip the colorbar below 0.5 and above 1.0 to enhance the visualization of plumes and downwellings. Furthermore, we always plot the non-dimensionalized temperature and the non-dimensionalized radius.

We find that ConvAE with a latent size of 840 (or width \times height \times channels = $5 \times 7 \times 24$) offers an excellent compression factor of 142, while being able to reconstruct the temperature fields with a mean relative accuracy of 99.80% on the test set. In comparison, ConvAEs with 1620- and 7600-dimensional latent spaces are 99.88% and 99.90% accurate, respectively. To calculate the mean relative accuracy, we dimensionalize the temperature using Eq. (A7) to avoid division by zeros.

In the subsequent sections, we discuss the training of ML algorithms to predict the temperature fields compressed to 840 numbers and then reconstruct the 302×394 -sized temperature fields using the trained decoder, i.e. we can compress the data from 1 TB to approximately 7 GB. We also tested the 1620- and 7600-dimensional encodings when predicting the compressed temperature fields - as explained in the subsequent section - and found no significant improvement in the predictions, suggesting that the accuracy of the 840-dimensional compression is good enough for a dataset of

this size. Keeping the latent space as small as possible is especially useful when training LSTMs, because an LSTM cell has 8 times as many trainable parameters as a dense FNN layer.

The results of the simple ConvAE are also encouraging in the light that the computational grid is structured, but not uniform, as is typically the case in computer vision applications. It seems that the different filters in the ConvAE are capable of capturing features at different spatial scales. For now, accounting for the curvilinear nature of the mesh and thereby potentially achieving higher compressibility and/or accuracy remains subject to future research.

IV. PREDICTION OF COMPRESSED TEMPERATURE FIELDS

A. Neural networks

With the temperature fields compressed, we now move on to predicting this latent space representation of 840 numbers, based on the five parameters of the simulation. Since [27] demonstrated that an FNN which takes the five parameters as inputs and time as a sixth input is capable of predicting the 1D temperature profile with a high accuracy, FNN is an obvious candidate for predicting the compressed temperature fields.

Fig. 5(a) shows the five parameters and the time that are used as inputs to the FNN to predict the compressed temperature fields. For computational efficiency, we optimize the network weights using small mini-batches of 16 temperature fields and Adam. Also, we used Scaled Exponential Linear Unit (SELU) as our activation function [58]:

$$SELU(x) = \lambda \begin{cases} x \\ \alpha e^x - \alpha \end{cases}, \quad (3)$$

as it seemed to deliver slightly better results than tanh. In Eq. (3), $\alpha = 1.67326324$ and $\lambda = 1.05070098$ are pre-defined and come from the original paper [58]. We further schedule the learning rate to decrease by a factor of 10 after 200 epochs and then again by a factor of 10 after the next 300 epochs. During the training, we only save the network weights if the validation loss drops. Furthermore, we employ a dropout of 5% after each hidden layer as a further regularization. Once the training has finished, the 2D fields are reconstructed from the predicted latent states as a post-processing step using the already trained decoder.

We trained different FNN architectures with fully connected dense layers. The evolution of the mean squared error (MSE) on the training and the cross-validation data is plotted in Fig. 6. The

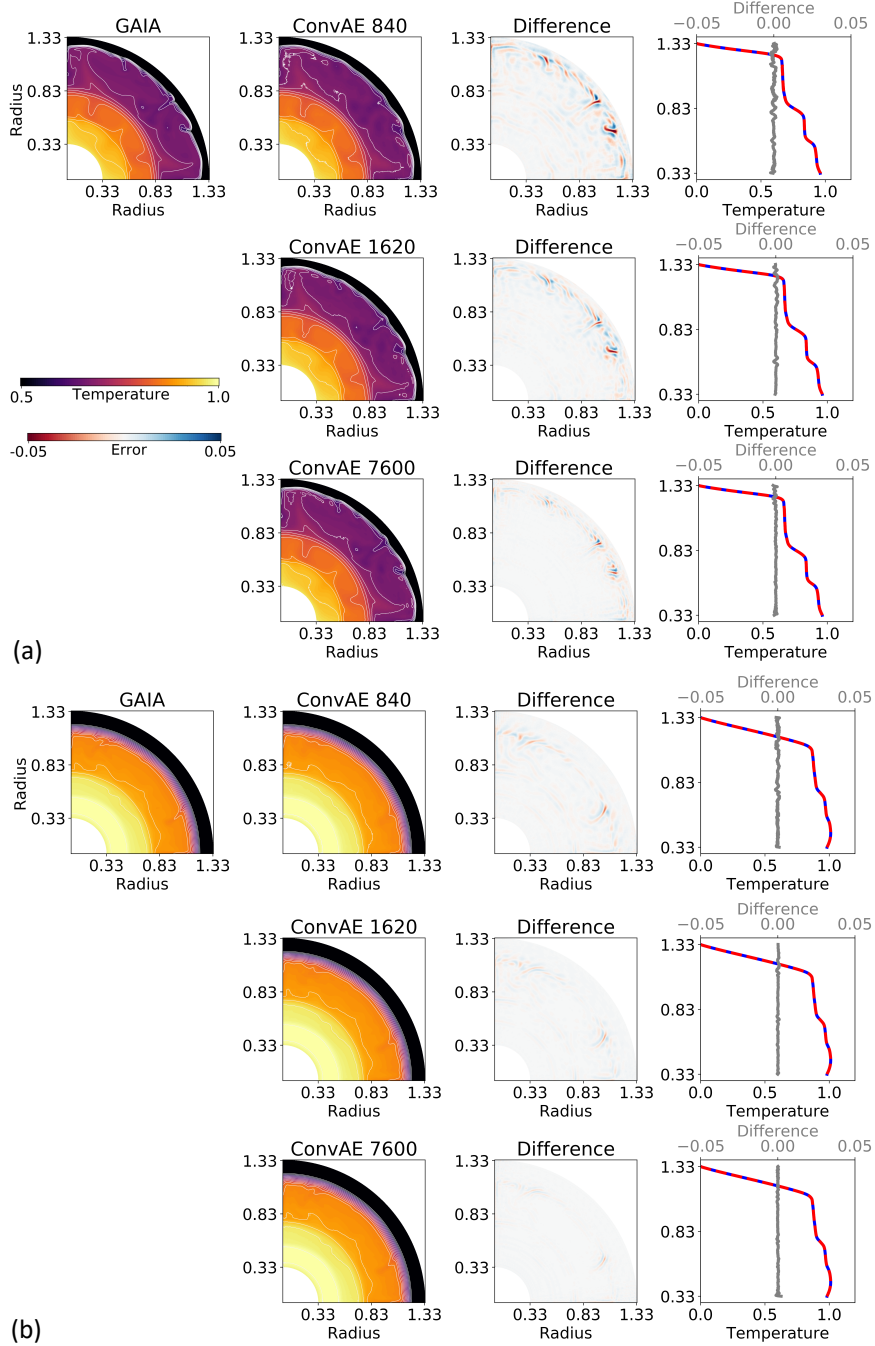


FIG. 4. (a) First example of compression and reconstruction from the test set for different architectures of convolutional autoencoders. ConvAE 840 represents a 142-fold compression of the original field (302×394), while ConvAE 1620 compresses the field by a factor of 73 and ConvAE 7600 by a factor of 16. (b) Second example with vigorous convection and small-scale structures from the test set. For each ConvAE, the error in the reconstruction is plotted in the third column, along with the horizontally-averaged 1D temperature profiles in the fourth column (and the difference between the two).

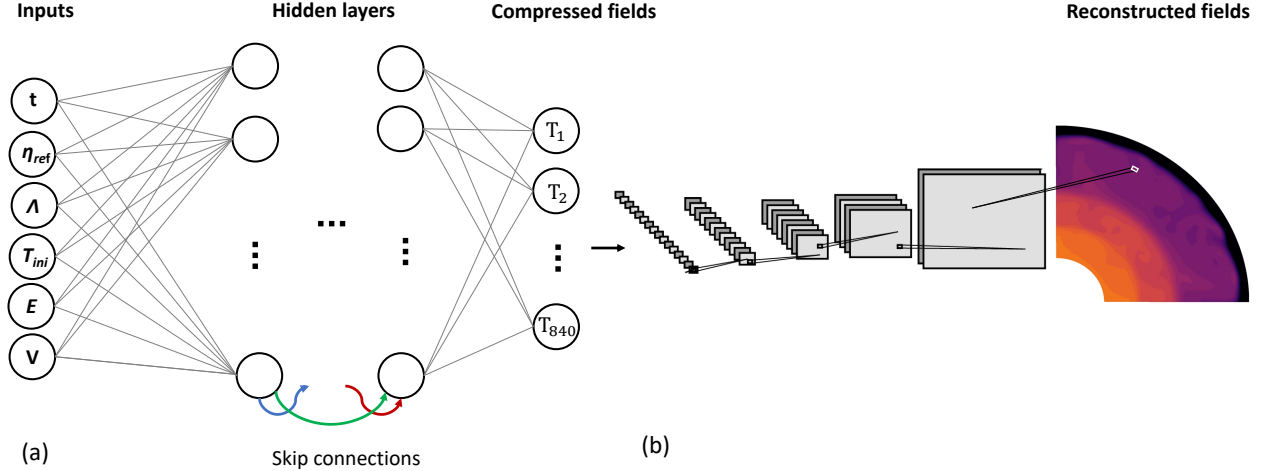


FIG. 5. (a) The five parameters governing mantle convection and time are used to predict the compressed temperature fields. Skip connections are used to add the output of each hidden layer after activation to that of each of the following hidden layers before activation. (b) After the training is complete, the trained decoder from ConvAE is used to reconstruct the temperature field back to its original dimensions.

cross-validation loss over epochs converges to approximately the same value, especially for deeper networks. Since, we also test some deep architectures such as five hidden layers with 800 neurons each and eight hidden layers with 400 neurons each, we use skip connections from each hidden layer to all the following hidden layers via addition. Given the challenging task of predicting a 840-dimensional vector, we use fairly deep architectures. Therefore, we employ SELU activation and skip connections to overcome the problem of vanishing gradients typically observed in deep networks.

For the FNN with 8 hidden layers of 400 units each, we take two examples from the test set with different convection patterns and plot one characterized by a sluggish behavior in Fig. 7 and another one characterized by vigorous convection in Fig. 8. A schematic of the this particular FNN architecture is provided in Appendix C. The figures show the 2D temperature fields computed with GAIA (first column), those predicted by the FNN (second column), the difference between the two (third column), as well as the horizontally-averaged 1D temperature profiles along with their difference (fourth column). While the 1D profiles show good agreement (as previously demonstrated by [27]), the 2D predictions show more inaccuracies. In particular, cold, sub-lithospheric downwellings, which are a fundamental feature of the planform of stagnant-lid convection, tend to be completely lost by the FNN prediction. See Supplemental Material at [59] for animations of these two examples along with three further examples from the test set.

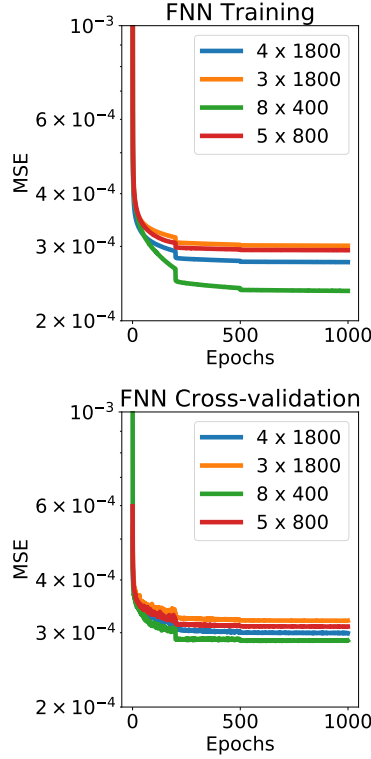


FIG. 6. The evolution of the mean-squared error (MSE) on (a) training data and on (b) cross-validation data for different FNN architectures. The legend shows the number of hidden layers as well as the number of neurons per layer of a given FNN architecture. For example, 4×1800 means the network has four hidden layers with 1800 neurons each. The step-like drop after 200 epochs is a result of the decrease in the learning rate.

In comparison with the GAIA simulations, the FNN predictions also fail to capture the vigor of convection. Even when the FNN captures a downwelling early on in the evolution, its lateral transport is not captured. On average, the 2D temperature fields predicted by FNN are 99.30% accurate (mean relative accuracy of dimensionalized temperature fields) with respect to GAIA and 99.35% with respect to ConvAE.

B. Long short-term memory (LSTM)

The failure of the FNNs to capture lateral convection structures such as down- and upwellings can be attributed in large part to the fact that the temporal snapshots of any given simulation are disconnected. By treating time as an additional input variable and shuffling time-steps of different simulations (but within the training/validation/test splits), the details of the dynamics of the flow

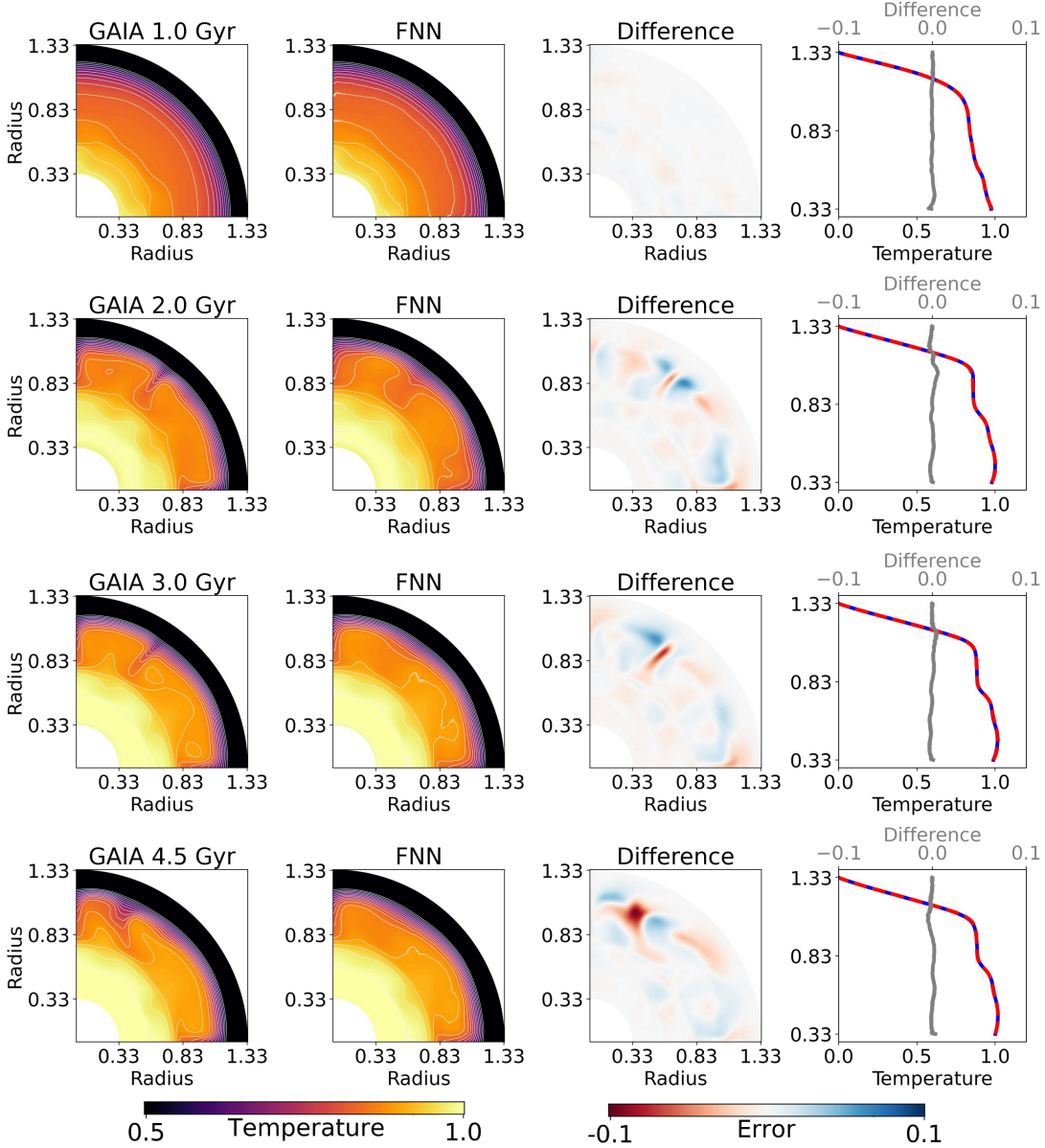


FIG. 7. Example of a sluggishly convecting mantle from the test set. The parameters are: $\eta_{\text{ref}} = 3.6 \times 10^{21}$ Pa S, $E = 1.6 \times 10^5$ J mol $^{-1}$, $V = 4.4 \times 10^{-6}$ m 3 mol $^{-1}$, $\Lambda = 15.3$ and $T_{\text{ini}} = 1634$ K. The temperature field from GAIA and its equivalent FNN prediction are shown in column 1 and 2, respectively. The third column shows the difference between the two. Column 4 shows the horizontally-averaged 1D temperature profiles from GAIA (solid blue) and FNN (dashed red), as well as the difference between the two (grey).

are largely lost.

Hence, we turn our attention to recurrent neural networks that have been shown to be successful for a variety of Natural Language Processing tasks. Recurrent architectures such as LSTM [60]

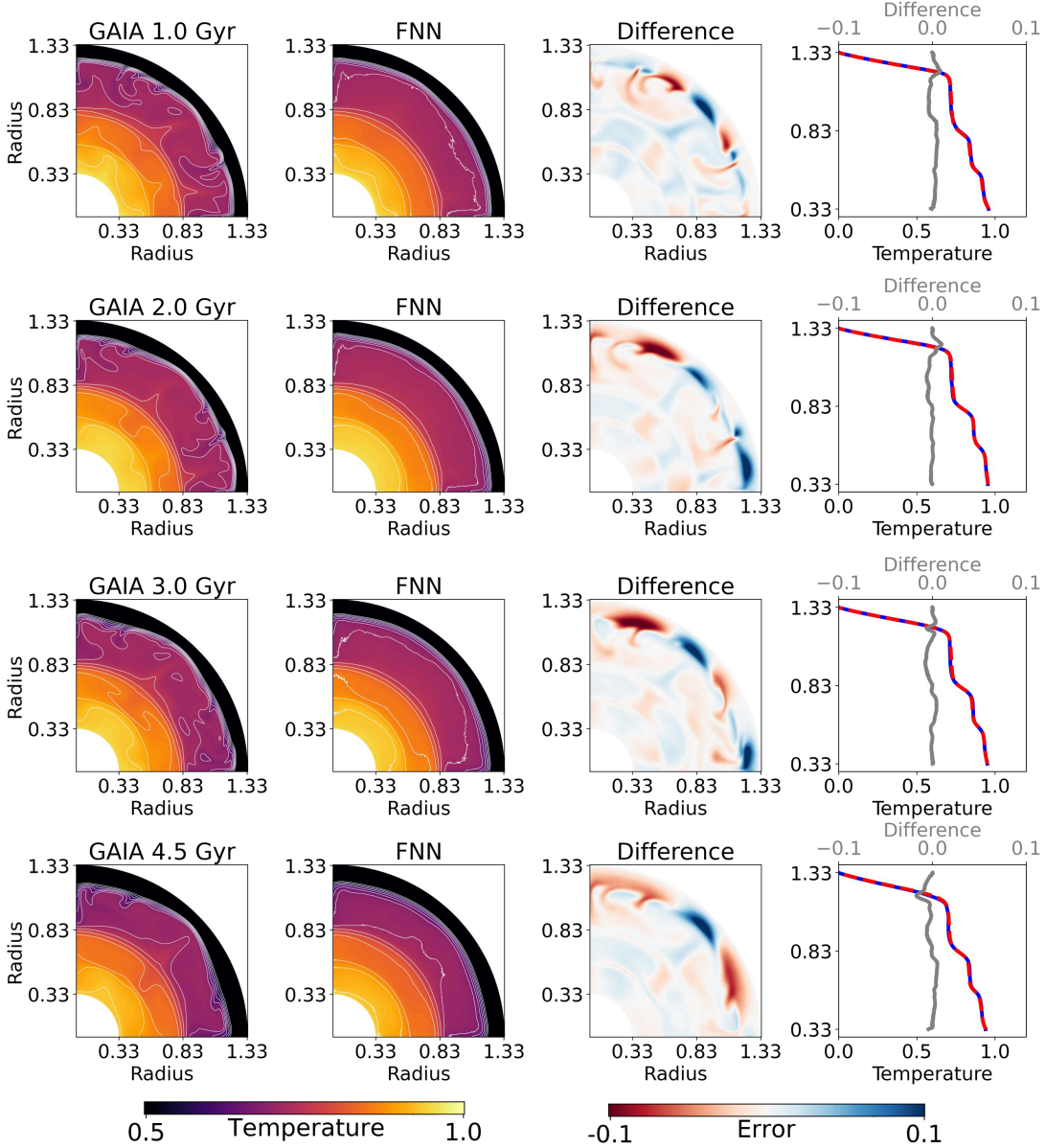


FIG. 8. Example of a vigorously convecting mantle from the test set. The parameters are: $\eta_{\text{ref}} = 5.0 \times 10^{19}$ Pa S, $E = 1.5 \times 10^5$ J mol $^{-1}$, $V = 7.6 \times 10^{-6}$ m 3 mol $^{-1}$, $\Lambda = 30.7$ and $T_{\text{ini}} = 1705$ K. The temperature field from GAIA and its equivalent FNN prediction are shown in column 1 and 2, respectively. The third column shows the difference between the two. Column 4 shows the horizontally-averaged 1D temperature profiles from GAIA (solid blue) and FNN (dashed red), as well as the difference between the two (grey).

provide a back-propagation mechanism acting through a sequence of inputs (such as time-steps of a simulation), thereby allowing the network to learn temporal dynamics (e.g., [34, 61]).

As shown in Fig. 9(b), LSTMs differ from an FNN architecture (Fig. 9(a)) in that they use a

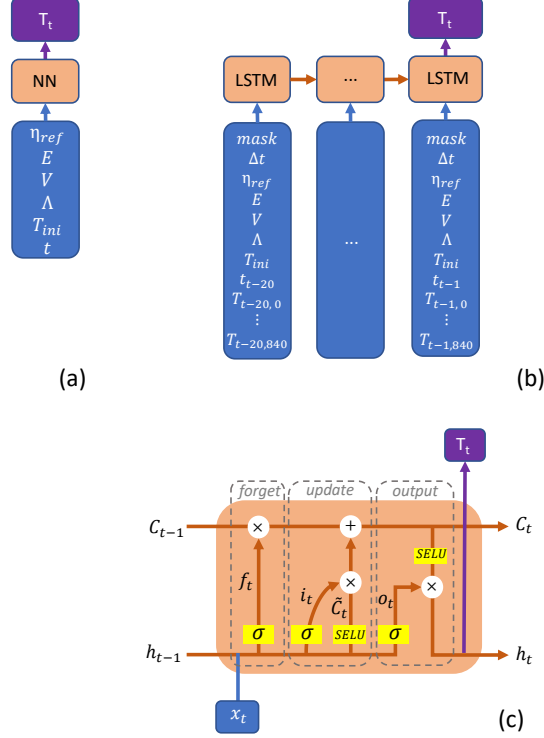


FIG. 9. Comparison of how (a) the FNNs are trained and how the (b) many-to-one LSTMs are trained for predicting compressed the temperature field T_t . (c) Each LSTM cell receives as input the compressed temperature field at time $t - 1$ along with other parameters: a mask for whether the next time-step exists, difference between the time-step used as input and the one being predicted (Δt), the five mantle convection parameters (η_{ref} , E , V , Λ and T_{ini}), as well as the time itself t_{t-1} of the input compressed temperature field. In practice, the mask is used by the data-generator to decide whether to provide the next time-step as output to train the LSTM on.

set of previous time-steps to predict the next time-step. Furthermore, each LSTM cell (Fig. 9(c)) is more elaborate than a simple neuron of a fully-connected FNN. [62] provides an accessible explanation of LSTMs. A brief description of an LSTM cell is provided in Sec. D.

As before, we use Keras API for setting up the forward graph (Eqs. (D1)–(D5)) and optimizing the trainable parameters by minimizing the MSE between the prediction and the target. Since, the time-steps for simulations were stored after a specified number of time-steps of the numerical solver, as well as at every physical time-interval of 100 million years of planetary evolution, we ended up with non-uniform time-series. Hence, we use “masking”, as done for example by [63], to specify if there is a time-step to predict ($mask = 1$) or not ($mask = 0$). Most of the simulations have less than 200 time-steps, although 9 simulations exceed 400 time-steps.

In addition to the mask, we also input the size of the time-step (Δt), the five model parameters, time and the compressed temperature field at the current time (t). After some trial and error, we found 20 time-steps to serve as a rich-enough input to predict the 21-st. When predicting a time-step when 20 previous inputs are not available, say for time-step 10, we took time-steps 0 through 9 and filled the rest with time-step 0. This way, the thermal evolution of a planet can be simulated based purely on input parameters. The input shape of each mini-batch is, thus, (simulations= 16, time-steps= 20, input= 848) and the output shape is (simulations= 16, time-steps= 1, output= 840).

The four different gates with two sets of matrices each in an LSTM cell mean that there are 8 times as many trainable parameters per hidden layer as a regular dense layer in an FNN. Therefore, we tested a limited number of LSTM architectures as shown in Fig. 10. For computational efficiency, we use *SELU* again, instead of *tanh* as activation function [64]. Such models take about 2 weeks on a Tesla V100 GPU to reach asymptotic loss values. Just as in the case of training FNNs, we use the following strategies to prevent over-fitting: (1) Storing the weights only if the validation loss drops, (2) using a dropout of 5% for each hidden layer, and (3) training on mini-batches of 16.

The difference between the loss curves for different architectures is small and given the stochasticity associated with training of LSTMs, not particularly enlightening with one small exception. The loss curve for the LSTM with five hidden layers of 900 cells each diverged around epoch 20, but seems to have found its way back a few epochs later, indicating that such a deep architecture with roughly 33 million trainable parameters might be prone to exploding gradients, especially when the learning rate is too large (0.0001). Of course, the initial learning rate is divided by a factor of 10 after the first 50 epochs and then by another factor of 10 after the next 150 epochs (if reached).

LSTM architectures achieve a lower MSE loss (Fig. 10) than FNN architectures (Fig. 6). However, it is important to note that the MSE loss displayed for LSTMs during training is based on a prediction, which takes the highly accurate compressed temperature fields as inputs. In inference mode, as the network iteratively predicts the next time-step based on the previously predicted 20 time-steps, this accuracy might decrease due to accumulation of error. Therefore, we calculated the mean relative accuracy for all the simulations in the cross-validation set using different LSTMs in purely inference mode and display them along with the mean relative accuracy on the test set in Table I. Table I seems to contradict the results of the MSE plots in Fig. 10. While during training, the [1800, 1800] LSTM with 46.5 million trainable parameters attained the lowest MSE,

TABLE I. mean relative accuracy of different LSTM architectures on the cross-validation (CV) and the test sets, when computed in inference mode (i.e. each input temperature field is also an LSTM prediction). For reference, mean relative accuracy for the FNN architecture is also presented.

Accuracy	wrt GAIA (%)	wrt ConvAE (%)
2 × 900		
Test	98.182 ± 9.853	98.237 ± 9.854
CV	96.542 ± 15.041	96.597 ± 15.049
3 × 900		
Test	99.222 ± 0.515	99.278 ± 0.513
CV	99.109 ± 0.602	99.164 ± 0.605
4 × 900		
Test	99.226 ± 0.524	99.285 ± 0.525
CV	99.082 ± 0.674	99.141 ± 0.682
5 × 900		
Test	99.199 ± 0.537	99.257 ± 0.537
CV	99.051 ± 0.680	99.108 ± 0.687
4 × 600		
Test	99.221 ± 0.495	99.281 ± 0.495
CV	99.081 ± 0.665	99.140 ± 0.671
2 × 1200		
Test	98.275 ± 8.313	98.328 ± 8.317
CV	97.364 ± 11.352	97.416 ± 11.356
2 × 1800		
Test	98.561 ± 6.446	98.616 ± 6.445
CV	98.326 ± 7.376	98.379 ± 7.380
FNN		
Test	99.297 ± 0.433	99.354 ± 0.422
CV	99.207 ± 0.482	99.262 ± 0.475

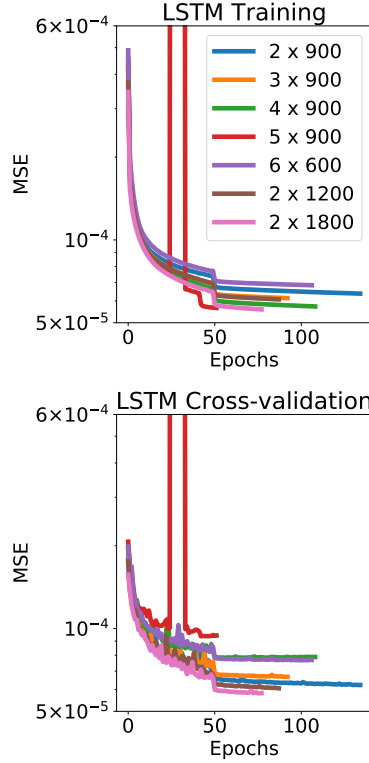


FIG. 10. The evolution of MSE on (top panel) training data and (bottom panel) cross-validation data for different LSTM architectures. The legend shows the number of hidden layers as well as the number of LSTM cells per layer. For example, 4×900 means the network has four hidden layers with 900 cells each. The step-like drop after 50 epochs is a result of the decrease in the learning rate.

in inference mode it had the third lowest mean relative accuracy of all architectures. In contrast, the $[600, 600, 600, 600]$ LSTM with 12.7 million trainable parameters, achieved a higher accuracy on the cross-validation set and on the test set. We simply do not have enough simulations to fit 46.5 million weights without over-fitting, which is evident from the high standard deviation of absolute relative accuracy on test and cross-validation sets. In inference mode, the errors in predicting time-steps can accumulate to the point where a simulation diverges. Luckily, this behavior is not observed in the $[600, 600, 600, 600]$ LSTM, for example. We present the evolution of the averaged MSE with time for different simulations in the test set, for both the FNN and LSTM, in Fig. 11. As the mantle typically starts cooling after some point in the thermal evolution, the convection should get slightly less vigorous. This means that the upwellings and downwellings should have a longer wavelength and therefore, become slightly easier to predict. Nevertheless, it seems that the lack of data towards the end of the evolution (from unfinished simulations) is the reason for the increased

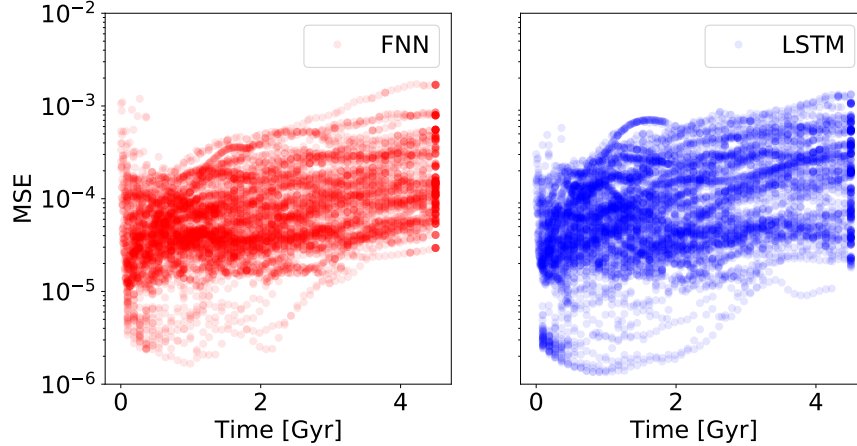


FIG. 11. Spatially averaged MSE for FNN (8×400) and LSTM (4×600) vs. physical time for each simulation in the test set.

error. In case of LSTMs it can be further exacerbated by the accumulation of error.

For both the LSTM and the FNN, the error tends to increase marginally with time. This can be attributed in large part to the fact that data get sparser with increasing time as not all simulations finished. For the LSTM though, the accumulation of error can be another contributing factor.

In Fig. 12 and Fig. 13, we plot the same two examples from the test set as in subsection IV A, but using the $[600, 600, 600, 600]$ LSTM this time. Despite the slightly lower mean relative accuracy compared to FNN, the LSTMs do a better job of capturing the convection structures. This is especially true for the more sluggish simulations such as the one in Fig. 12, where the downwelling is not only formed at 2 Gyr (second row), but also maintained and transported towards the left boundary in time, unlike the FNN.

Fig. 13 is a good example of how the small-scale downwellings formed under more vigorous convection are not well captured by the LSTM. This explains the richness of structures one can see from the difference plots in column 3. Nevertheless, the big downwelling captured at 1 Gyr to the right of the domain at radius of 0.83 to 1.1, for example, or the upwelling at same radial location, but towards the middle of the domain present an improvement over the smudged-out prediction of the same simulation by an FNN in Fig. 8.

In summary, LSTMs are better at predicting sharper structures such as downwellings as well as the dynamics of their transport. A reason for the lower relative mean accuracy compared to FNNs is that the movement of plumes and downwellings, while captured, can be longitudinally off. In other words, the downwelling captured in Fig. 12 (rows 2 and 3) are slightly shifted in the

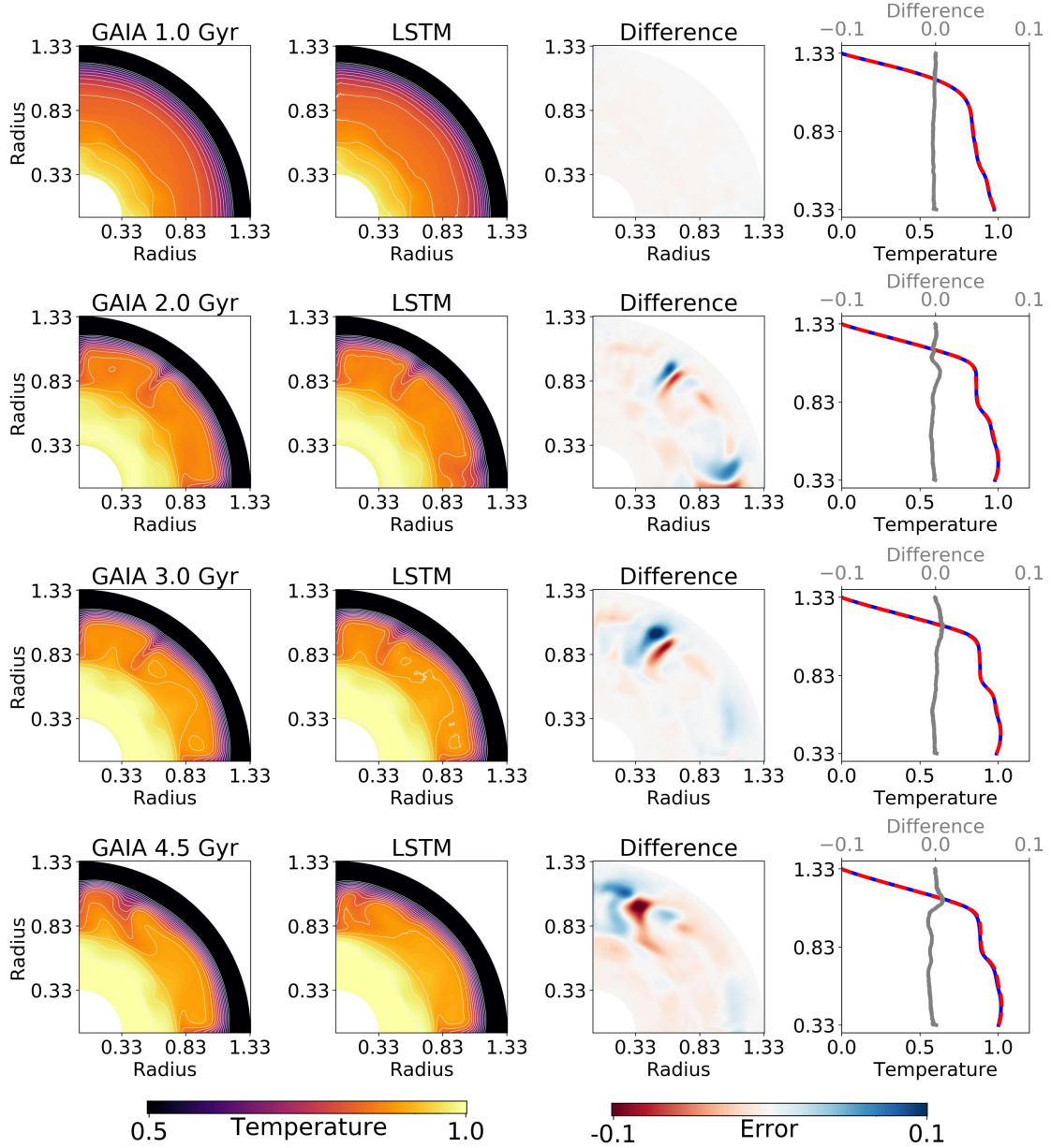


FIG. 12. Example 1 from the test set. The temperature field from GAIA and its equivalent LSTM prediction are shown in column 1 and 2, respectively. The third column shows the difference between the two. Column 4 shows the horizontally-averaged 1D temperature profiles from GAIA (solid blue) and LSTM (dashed red), as well as the difference between the two (grey).

angular direction. The same can be seen in all the difference plots of Fig. 13, at radial locations of 0.5 – 0.83 for a downwelling and 0.83 – 1.2 for a plume in the longitudinal center.

One could further examine the 1D temperature profiles, obtained by horizontally averaging the 2D temperature fields (column 4 in Fig. 7, Fig. 8, Fig. 12, Fig. 13). LSTM temperature profiles

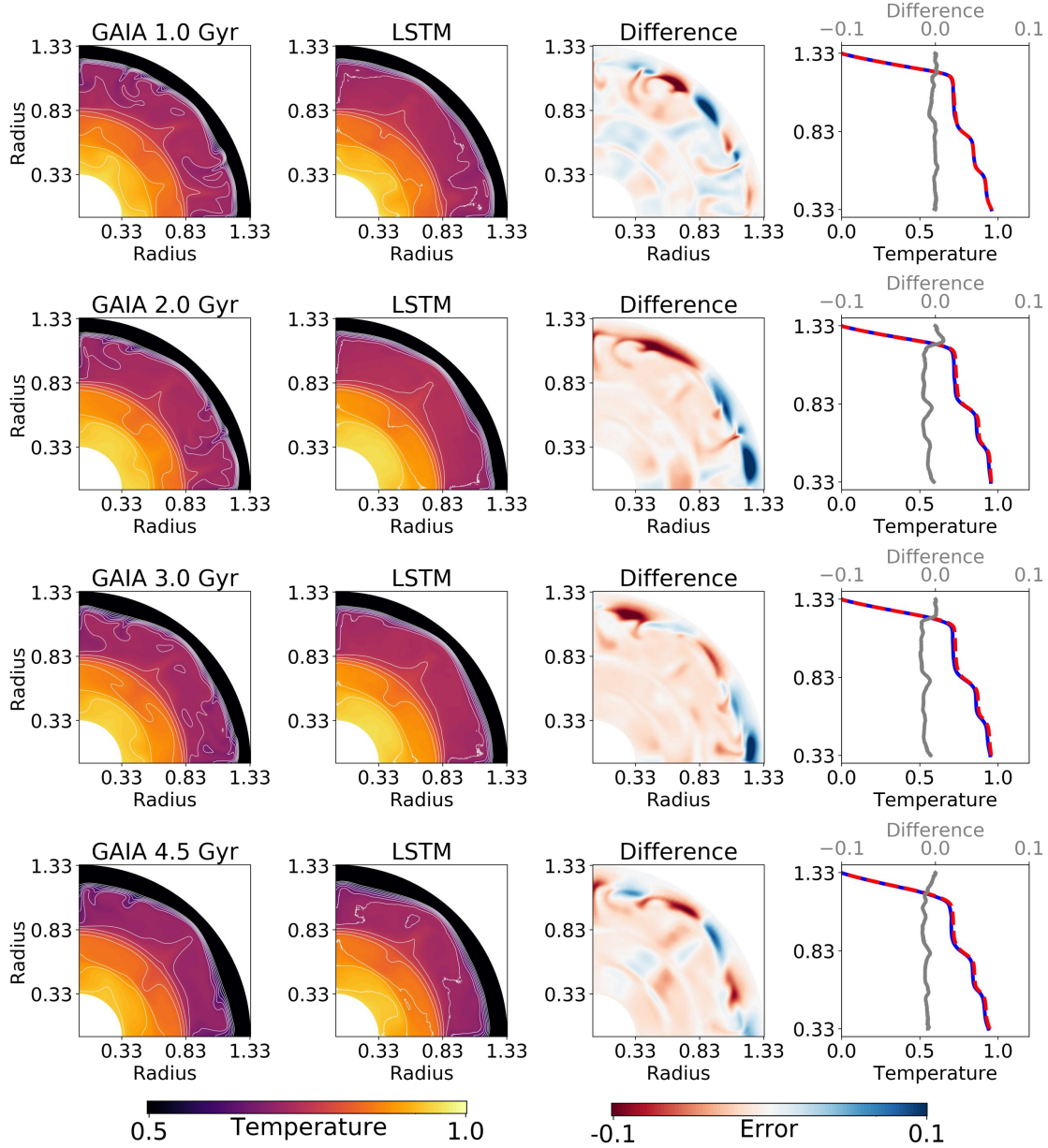


FIG. 13. Example 2 from the test set. The temperature field from GAIA and its equivalent LSTM prediction are shown in column 1 and 2, respectively. The third column shows the difference between the two. Column 4 shows the horizontally-averaged 1D temperature profiles from GAIA (solid blue) and LSTM (dashed red), as well as the difference between the two (grey).

have a mean relative absolute accuracy of 99.42%, while those of FNN are 99.71%. While this intuitively counters the longitudinal shift argument, the mean of smudged-out/lacking plumes and downwellings predicted by the FNN can match the mean of sharp plumes and downwellings of GAIA better than the LSTM. The fact that the LSTM temperature fields capture some, but not all the

downwellings in cases of vigorous convection (Fig. 13) can throw the horizontal mean off. In fact, finding an error metric that is invariant to, for example, longitudinal shift of a downwelling is non-trivial. [65], for example, show how modified versions of normalized root-mean-squared (NRMSE) can be computed that are invariant to certain effects such as multiplication by a constant, or phase shift for image reconstruction. Since, from a planetary evolution viewpoint, the magnitude of the temperature field matters, one could attempt to find similar shift- or rotation-invariant metrics but in terms of MSE, instead of NRMSE. Furthermore, one must consider whether to use the modified MSE expression to only evaluate the error, or also to optimize the weights of the machine learning architectures. In the next subsection, we show that MSE seems capable of learning some non-trivial dynamics of mantle convection, as long as the underlying machine learning algorithm is suitable.

The mean relative accuracy for all the simulations in the test set is provided in Fig. 14. Low accuracy seems to be slightly correlated to low reference viscosity and low activation energy for the diffusion creep. A low reference viscosity generally leads to more vigorous convection, thereby inducing small-scale convection structures, which the LSTM finds difficult to predict. Similarly, a low activation energy, or equivalently, a low dependence of viscosity on temperature has the same qualitative effect of reducing viscosity. Otherwise, the method works well across the entire range of parameters.

C. POD comparison of FNN and LSTM predictions

We analyze the FNN and LSTM predictions from the lens of POD (proper orthogonal decomposition). Following e.g. [66], we compute the Singular Value Decomposition of a “tall” simulation matrix $X \in \mathbb{R}^{p \times q}$ (spatial points \times time-steps):

$$X = U\Sigma V^*, \quad (4)$$

to obtain the spatial modes $U \in \mathbb{R}^{p \times r}$, complex conjugate V^* of temporal modes $V \in \mathbb{R}^{r \times q}$ and the POD coefficients (eigenvalues) $\Sigma \in \mathbb{R}^{r \times r}$, where r is determined by the minimum of p and q . In Fig. 15(a), we plot the eigenvalues for the simulation in Fig. 12 and Fig. 7 and in Fig. 15(b), we display the eigenvalues for the simulation in Fig. 13 and Fig. 8.

As can be seen in both cases in Fig. 15, the eigenvalues of an FNN-predicted temperature field evolution decay very rapidly after the first three to five modes. Hence, the cumulative distribution function (CDF) of the FNN predictions is the steepest, reaching most of its energy within the first

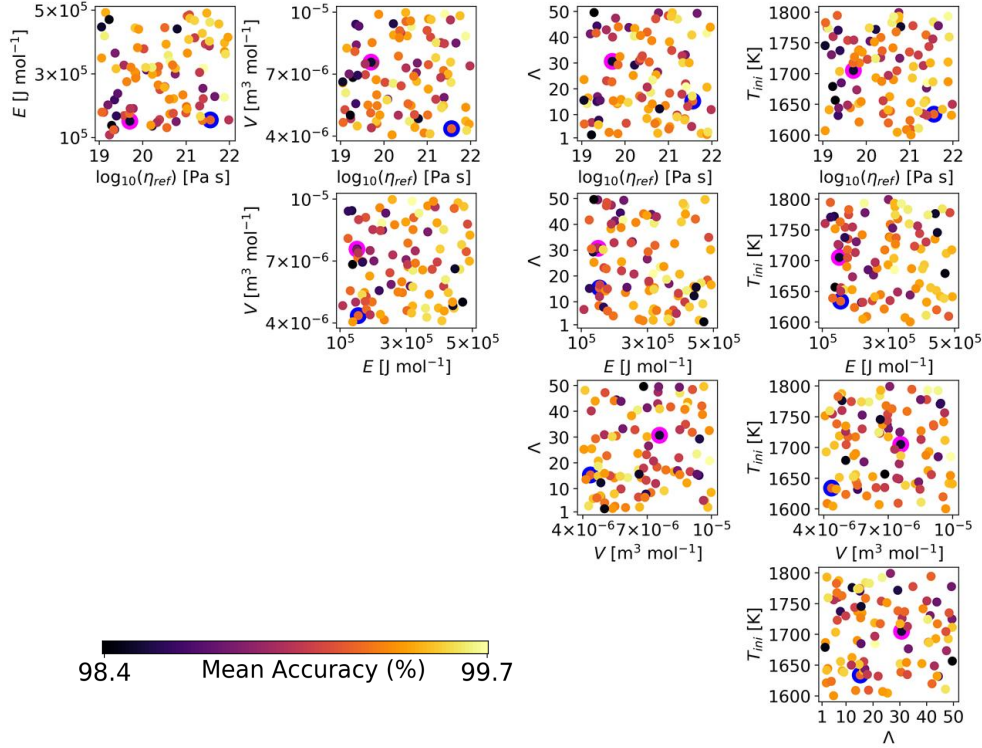


FIG. 14. Mean relative accuracy of LSTM predictions of the temperature fields for all the simulations in the test set with respect to the original GAIA simulations, expressed as a percentage. The mean relative accuracy is plotted with respect to two parameters at a time - indicated on the x- and y-axis with the units in which the parameters are measured. Example in Fig. 12 is circled in blue, while the example in Fig. 13 is circled in magenta.

few modes, as opposed to the other CDFs where latter modes also carry non-negligible energy. This phenomenon can be visually observed in the animations of all the five test-case simulations. See Supplemental Material at [59] for the animations. The FNN predictions are not “energetic” enough. The POD coefficients of the simulations predicted by LSTM decay less rapidly, even when, occasionally, the decay is desired (Fig. 15(a) modes 40–76). However, in the case of vigorous convection, Fig. 15(b) shows that LSTMs, while better than FNNs, still do not fully capture the energy characteristics of the GAIA simulation. On average, the sum of eigenvalues of the FNN predictions on the entire test set amounts to 96.51% of the the sum of eigenvalues of the GAIA predictions. For LSTM-predicted simulations, the sum of POD coefficients increases to 97.66% relative to those of GAIA simulations. Thus, LSTMs capture the dynamics of the simulations better, while FNNs provide a marginally better mean accuracy for snapshots.

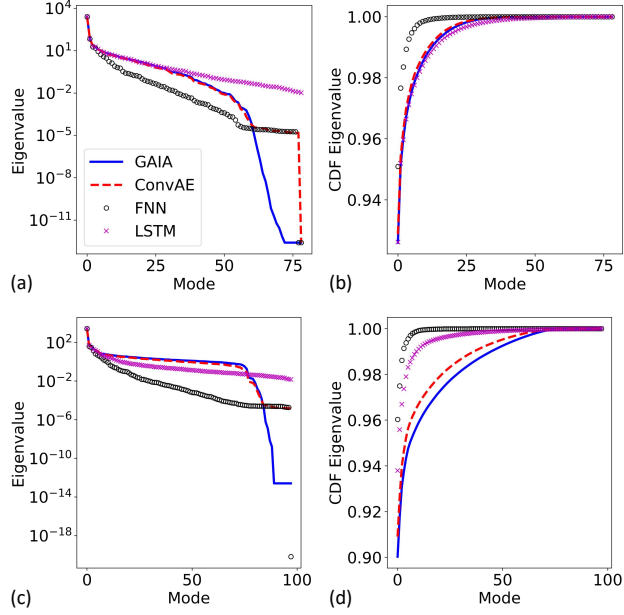


FIG. 15. (a) POD coefficients and (b) their cumulative distribution for example simulation 1 in the test set and (c) and (d) correspond to example simulation 2.

V. CONCLUSION AND FUTURE WORK

We use deep learning techniques to model parameterized surrogates of mantle convection simulations. The data to train the algorithms comes from 10,525 mantle convection simulations of a Mars-like planet, run on a 2D quarter-cylindrical grid. Focusing on only one state variable in this study - temperature - we first compressed the temperature fields using convolutional autoencoders by a factor of 142 (from 1 TB to 7 GB) and then tested two regression algorithms for predicting the compressed temperature fields given five key parameters: reference viscosity (linked to the Rayleigh number), activation energy and activation volume of the diffusion creep, an enrichment factor for radiogenic elements in the crust and the initial mantle temperature (see Fig. 1(a)).

We found that while feedforward neural networks (FNN) offer a reasonable mean accuracy, the sharper plumes and downwellings formed in the upper mantle are never transported laterally and are often lost early in the evolution. In contrast, many-to-one long-short term memory networks (LSTM) were able to capture the sharper convection structures along with their lateral transport more often, but ultimately, delivered a slightly lower mean accuracy (99.22%) in comparison to FNNs (99.30%). Two factors are mainly responsible for the lower mean accuracy: (1) the convection structures were longitudinally shifted and (2) the prediction error can accumulate in

time. Despite this, the eigenvalues obtained through proper orthogonal decomposition (POD) show that the FNN predictions decay too rapidly after the first three to five modes, while LSTM predictions decay less rapidly and hence capture the flow dynamics more accurately than the FNN, if not perfectly. When summed, the eigenvalues from FNN predictions and the eigenvalues from LSTM predictions amount to 96.51% and 97.66% relative to those obtained through POD of the original simulations, respectively.

This study serves as a first-proof that deep learning can be used to model high-dimensional parameterized surrogates of mantle convection simulations. Given five parameters, the complete spatio-temporal evolution of the temperature field can be predicted up to a reasonable accuracy, i.e. the longer wavelength structures such as the 1D temperature profile and larger plumes and downwellings as well as their lateral transport can be captured, albeit not perfectly. With respect to the thermal evolution of terrestrial planets, the 2D temperature fields can be used to calculate a number of fields of interest and relate them to various quantities that can be inferred from actual observations. Lateral variations in the heat flux are important for estimating the elastic lithospheric thickness (e.g. [52]). Spatio-temporal variations in the heat flux at the core-mantle boundary affect the generation and morphology of the magnetic field (e.g. [67]). The formation of plumes and downwellings is important for calculating the amount of melt produced during the thermal evolution and to relate this to estimates of the thickness of the crust (e.g. [68]). The fact that plumes and downwellings are not accurately captured by the FNN can impact local melt production. Similarly, LSTM's longitudinally shifted plumes compared to the true simulations can result in slightly different looking crustal distributions. It is not straightforward to predict how consequential these errors would be to constraining the parameters. Ideally, one would conduct an inverse study to test the sensitivity of uncertainties in the observables resulting from instrumentation and/or from the surrogate model (e.g. [13]).

Parameterized surrogates such as the ones presented in this paper or the one proposed in [37] are primarily useful for performing parameter-studies - be it placing constraints on the evolution of a planet like Mars, or optimizing an airfoil to achieve the target aerodynamic performance.

This is fundamentally different from applications where the time-steps of the same simulation can be split into training and test sets (e.g., [32–34]). On the one hand, the latter is an easier learning task because one can expect the dynamics of a single simulation to exist on a significantly smaller manifold than multiple simulations with a wide range of parameters. On the other hand, the simplicity of the flow in our simulations (e.g. no turbulence, limited compressibility and 2D

flow instead of 3D) begs the question if parameterized surrogates can be learned for more complex flows than the one used in this study.

Particularly important would be the computational cost of generating the data. For flows similar to those considered here, but in 3D, running 10,000 simulations would be intractable. Worse yet, 10,000 simulations could be an order of magnitude less than needed to learn spatio-temporal dynamics in 3D. This is especially important in the light of the fact that while the 2D convection models provide significantly more information than 0D (e.g., [18–25]) or 1D evolution models [27], they still cannot be used to constrain parameters based on localized observational constraints in 3D such as crustal thickness, elastic lithospheric thickness or surface heat flux. Hence, this study only serves as a stepping-stone to the ultimate goal of using high-dimensional forward surrogates for probabilistic inversion of mantle convection parameters. To reach this goal, further research is needed into data-efficient methods for parameterized surrogate modelling of 3D mantle convection. One or both of the following approaches might hold the key. (1) One could use the partial differential equations and boundary conditions governing mantle convection to place soft and/or hard constraints on the optimization problem at hand (e.g., [33, 69–71]). (2) Although POD bases are not appropriate for this problem because they do not generalize well among simulations with different parameters, the idea of finding a set of basis functions, that one only needs to learn the coefficients to, remains an attractive one (e.g., [72–74]).

Despite these challenges, the combination of scientific deep learning and modern super-computing hardware presents an unprecedented opportunity for gaining robust insights into geophysical flows.

A JupyterNotebook to predict the entire spatio-temporal evolution of the 2D temperature field from five parameters is available on Github: https://github.com/agsiddhant/ForwardSurrogate_Mars_2D [75].

ACKNOWLEDGMENTS

We would like to thank the two anonymous reviewers whose comments helped improve a previous version of this paper. We acknowledge the support of the Helmholtz Einstein International Berlin Research School in Data Science (HEIBRiDS). We also acknowledge the North-German Supercomputing Alliance (HLRN) for providing HPC resources (project id: bep00087). This work was also funded by the German Ministry for Education and Research as BIFOLD - Berlin

Institute for the Foundations of Learning and Data (ref. 01IS18025A and ref 01IS18037A). We thank Klaus-Robert Müller for the useful discussion on methods used in this paper.

We list the author contributions following the taxonomy by [76]. *Conceptualization*: N.T., D.B., G.M.; *Methodology*: S.A., P.K., N.T.; *Software*: S.A.; *Validation*: S.A.; *Investigation*: S.A.; *Data curation*: S.A.; *Writing–Original Draft*: S.A., N.T.; *Writing–Review & Editing*: S.A., N.T., P.K., D.B., G.M.; *Visualization*: S.A., N.T.; *Supervision*: N.T., D.B., P.K., G.M.; *Funding Acquisition*: N.T., D.B., G.M.

Appendix A: Mathematical model of mantle convection simulations

The mathematical model used to run the simulations is the same as the one presented in [27]. Yet, for completeness, we describe it in detail in this appendix. Numerical values of the model parameters that are shared by all simulations are listed at the end of the appendix in Table II and Table III.

1. Governing equations

The equations of conservation of mass, linear momentum and thermal energy under the extended Boussinesq approximation can be written in non-dimensional form as follows (e.g. [40, 77]):

$$\nabla' \cdot \mathbf{u}' = 0, \quad (\text{A1})$$

$$\begin{aligned} -\nabla' p' + \nabla' \cdot \left[\eta' \left(\nabla' \mathbf{u}' + (\nabla' \mathbf{u}')^T \right) \right] \\ + \left(Ra \alpha' T' - \sum_{l=1}^3 Rb_l \Gamma_l \right) \mathbf{e}_r = 0, \end{aligned} \quad (\text{A2})$$

$$\begin{aligned} \frac{DT'}{Dt'} - \nabla' \cdot (k' \nabla' T') - Di \alpha' (T' + T'_0) u'_r - \frac{Di}{Ra} \Phi' \\ - \sum_{l=1}^3 Di \frac{Rb_l}{Ra} \frac{D\Gamma_l}{Dt} \gamma_l (T' + T'_0) - \frac{Ra_Q}{Ra} = 0, \end{aligned} \quad (\text{A3})$$

Primed quantities are non-dimensional. \mathbf{u}' is the velocity vector, p' the dynamic pressure, η' the viscosity, Ra the thermal Rayleigh number, α' the thermal expansivity, T' the temperature, Rb_l the Rayleigh number associated with the l -th phase transition, Γ_l the corresponding phase function (e.g. [44]), \mathbf{e}_r the unit vector in the radial direction, t' the time, k' the thermal conductivity, Di

the dissipation number, T'_0 the surface temperature, u'_r the radial component of the velocity, Φ' the viscous dissipation, and Ra_Q the Rayleigh number for internal heating. Viscosity, thermal expansivity and thermal conductivity are function of pressure and temperature (see eq. (1), Sec. A 4).

2. Non-dimensionalization of state variables

In eqs. (A1)–(A3), the dimensional variables are scaled as follows:

$$\mathbf{u}' = \mathbf{u} \frac{\rho_m c_{p_m} D}{k_{\text{ref}}}, \quad (\text{A4})$$

$$p' = p \frac{\rho_m c_{p_m} D^2}{\eta_{\text{ref}} k_{\text{ref}}}, \quad (\text{A5})$$

$$t' = t \frac{k_{\text{ref}}}{\rho_m c_{p_m} D^2}, \quad (\text{A6})$$

$$T' = \frac{T - T_0}{\Delta T}. \quad (\text{A7})$$

In eqs. (A4)–(A7), ρ_m is the mantle density and c_{p_m} its heat capacity; $D \equiv R_p - R_c$ is the mantle thickness (R_p and R_c are the planet and core radius, respectively); k_{ref} is the reference thermal conductivity, η_{ref} the reference viscosity (eq. 1) and ΔT the initial temperature drop across the mantle.

3. Dimensionless quantities

The following non-dimensional numbers appear in eqs. (A1)–(A3):

$$Ra = \frac{\rho_m^2 c_{p_m} \alpha_{\text{ref}} g \Delta T D^3}{\eta_{\text{ref}} k_{\text{ref}}}, \quad (\text{A8})$$

$$Rb_l = \frac{\rho_m c_{p_m} \Delta \rho_l g D^3}{\eta_{\text{ref}} k_{\text{ref}}}, \quad (\text{A9})$$

$$Ra_Q = \frac{\rho_m^3 c_{p_m} \alpha_{\text{ref}} g H_0 D^5}{\eta_{\text{ref}} k_{\text{ref}}^2}, \quad (\text{A10})$$

and

$$Di = \frac{\alpha_{\text{ref}} g D}{c_{p_m}}, \quad (\text{A11})$$

where, α_{ref} and k_{ref} are the reference thermal expansivity and conductivity, g is the gravitational acceleration, $\Delta\rho_l$ is the density contrast across the l -th phase-transition, and H_0 is the initial rate of mantle heat production due to radiogenic elements.

4. Thermal expansion and conductivity

The temperature- and pressure-dependent thermal expansivity and conductivity are calculated using the parametrizations introduced by [45], which in dimensional form, read:

$$\alpha(T, P) = \left(a_0 + a_1 T + a_2 T^{-2} \right) \exp(-a_3 P), \quad (\text{A12})$$

$$k(T, P) = (c_0 + c_1 P) \left(\frac{300}{T} \right)^{c_2}. \quad (\text{A13})$$

Here, a_0, \dots, a_3 and c_0, \dots, c_2 are coefficients based on experimental data valid for Mg-rich olivine. Numerical values of these coefficients are listed in Table III.

5. Phase transitions

We consider two solid-solid phase transitions in the olivine system, α to β -spinel and β to γ -spinel [44]. The temperature-dependent depth of the l -th phase boundaries $z_l(T)$ is calculated as:

$$z_l(T) = z_l^0 + \gamma_l (T - T_l^0). \quad (\text{A14})$$

Here, γ_l is the Clapeyron slope (positive for both phase transitions), z_l^0 is the reference transition depth and T_l^0 the corresponding reference temperature. The phase-transition function Γ_l is expressed in terms of z_l and phase transition width d_l as:

$$\Gamma_l = \frac{1}{2} \left(1 + \tanh \left(\frac{z - z_l(T)}{d_l} \right) \right). \quad (\text{A15})$$

6. Depletion of heat-producing elements and partial melting

We assume that a crust of thickness d_{cr} formed early in the planet evolution [41] and that this event led to the extraction of a large amount of radiogenic elements from the mantle [30]. We use an enrichment factor Λ to modify the bulk abundance of heat-producing elements C_0 in the mantle

(based on [42]) to a new depleted composition C_{depleted} as follows:

$$C_{\text{depleted}} = \frac{M_m C_0}{M_{\text{cr}} (\Lambda - 1) + M_m}, \quad (\text{A16})$$

where M_m and M_{cr} are the mass of the mantle and crust, respectively.

We model partial melting following the approach of [43]. Whenever the mantle temperature T locally exceeds the solidus, we calculate the melt fraction φ_i by solving the following equation:

$$c_p (T_i - T_{\text{sol}}) = L_m \varphi_i + c_p \Delta \varphi_i \Delta T_{\text{liq-sol}} (1 - \varphi_i), \quad (\text{A17})$$

where, T_i is the temperature in the i -th cell of the computational domain, T_{sol} the local solidus temperature, L_m the latent heat of melting and $\Delta T_{\text{liq-sol}}$ the local difference between the liquidus and solidus temperature. For the solidus and liquidus, we use the parameterization of [78] and of [79], respectively:

$$T_{\text{sol}} = e_0 + e_1 P + e_2 P^2 + e_3 P^3 + e_4 P^4, \quad (\text{A18})$$

$$T_{\text{liq}} = f_0 + f_1 P + f_2 P^2 + f_3 P^3 + f_4 P^4, \quad (\text{A19})$$

where T_{sol} and T_{liq} are the dimensional solidus and liquidus temperatures, respectively, P is the dimensional hydrostatic pressure and e_0, \dots, e_4 and f_0, \dots, f_4 are numerical coefficients listed in Table III.

Using the sum of melt produced in all cells at time-step t , φ_t , we adjust the internal heating Rayleigh number to model the further extraction of heat-producing elements due to melting:

$$Ra_{Q_t} = Ra_{Q_{t-1}} (1 - \Lambda \varphi_t), \quad (\text{A20})$$

7. Evolution of the core-mantle boundary temperature

An isothermal boundary condition is imposed at the core-mantle boundary whose temperature T_c evolves according to (e.g. [18]):

$$c_{p_c} \rho_c V_c \frac{dT_c}{dt} = -q_c A_c. \quad (\text{A21})$$

Here c_{p_c} is the specific heat-capacity of the core, V_c the volume of the core, q_c the average heat flux at the core-mantle boundary (CMB), and A_c the outer area of the core.

8. Rescaling of the core

As often done in simulations of mantle convection carried out in a 2D cylindrical shell geometry, a geometric rescaling of the core radius is applied in order to better reproduce the temperature field that would be obtained in a 3D spherical shell. Following [80], the radius of the core of the cylinder (R_c^{cyl}) is re-scaled in such a way that the core-to-planet radius ratio is the same as the core-to-planet surface ratio of a sphere. In formulas:

$$\left(\frac{R_c}{R_p}\right)^2 = \frac{R_c^{\text{cyl}}}{R_p^{\text{cyl}}} \quad (\text{A22})$$
$$R_p^{\text{cyl}} + R_c^{\text{cyl}} = 1,$$

where, R_p , R_c and R_p^{cyl} are respectively, the radii of the spherical planet, spherical core and cylindrical planet.

Appendix B: Distribution of parameters in the dataset

Fig. 16 plots the distribution of simulation parameters in the training, cross-validation and test sets.

Appendix C: FNN schematic

Fig. 17 is a schematic of one of the FNN architectures that were trained and the one we present the results for in Fig. 7 and Fig. 8. The trained FNN has 8 hidden layers with 400 neurons each and each hidden layer is connected to all the following hidden layers via skip connections. For example, hidden layer 0 is connected to hidden layers 1 through 7. For ease of visualization, we show only 3 hidden layers in Fig. 17.

Appendix D: Equations of an LSTM cell

Referring to Fig. 9(c), an LSTM cell has three main blocs. The “forget” gate f_t determines how much information from the previous cell-state C_{t-1} should be retained given the input vector x_t and the previous hidden state h_{t-1} (or equivalently, the previous output of an LSTM cell):

$$f_t = \sigma (W_f x_t + U_f h_{t-1} + b_f), \quad (\text{D1})$$

TABLE II. Values of fixed parameters shared by all simulations (part 1).

Parameter	Physical meaning	Value	Unit
$\Delta T_{t=0}$	¹ Initial temperature difference between core and surface	2000	K
T_0	¹ Surface temperature	250	K
ρ_c	¹ Core density	7000	kg m ⁻³
ρ_m	¹ Mantle density	3500	kg m ⁻³
c_{p_c}	¹ Core specific heat capacity	850	J kg ⁻¹ K ⁻¹
c_{p_m}	¹ Mantle specific heat capacity	1200	J kg ⁻¹ K ⁻¹
k_{ref}	¹ Reference thermal conductivity	4	W m ⁻¹ K ⁻¹
α_{ref}	¹ Reference thermal expansivity	2.5×10^{-5}	K ⁻¹
R_c	¹ Outer radius of the core	1700	km
R_p	¹ Planetary radius	3400	km
d_{cr}	Thickness of the crust	64.3	km
z_{ref}	Reference depth for viscosity	232	km
T_{ref}	¹ Reference temperature for viscosity	1600	K
$z_{\alpha\beta}^0$	¹ Reference depth for α to β spinel	1020	km
$z_{\beta\gamma}^0$	¹ Reference depth for β to γ spinel	1360	km
$\Delta\rho_{\alpha\beta}^0$	¹ Density difference for α to β spinel	250	kg m ⁻³
$\Delta\rho_{\beta\gamma}^0$	¹ Density difference for β to γ spinel	150	kg m ⁻³
$\gamma_{\alpha\beta}$	¹ Clapeyron slope for α to β spinel	3×10^6	Pa
$\gamma_{\beta\gamma}$	¹ Clapeyron slope for β to γ spinel	5.1×10^6	Pa
$T_{\alpha\beta}$	¹ Reference temperature for α to β spinel	1820	K
$T_{\beta\gamma}$	¹ Reference temperature for β to γ spinel	1900	K
d_l	¹ Width of phase transitions	20	km
${}^U C_0$	² Bulk abundance of uranium	16×10^{-9}	kg kg ⁻¹
${}^{\text{Th}} C_0$	² Bulk abundance of thorium	56×10^{-9}	kg kg ⁻¹
${}^{\text{K}} C_0$	² Bulk abundance of potassium	305×10^{-6}	kg kg ⁻¹

¹[29] ²[81].

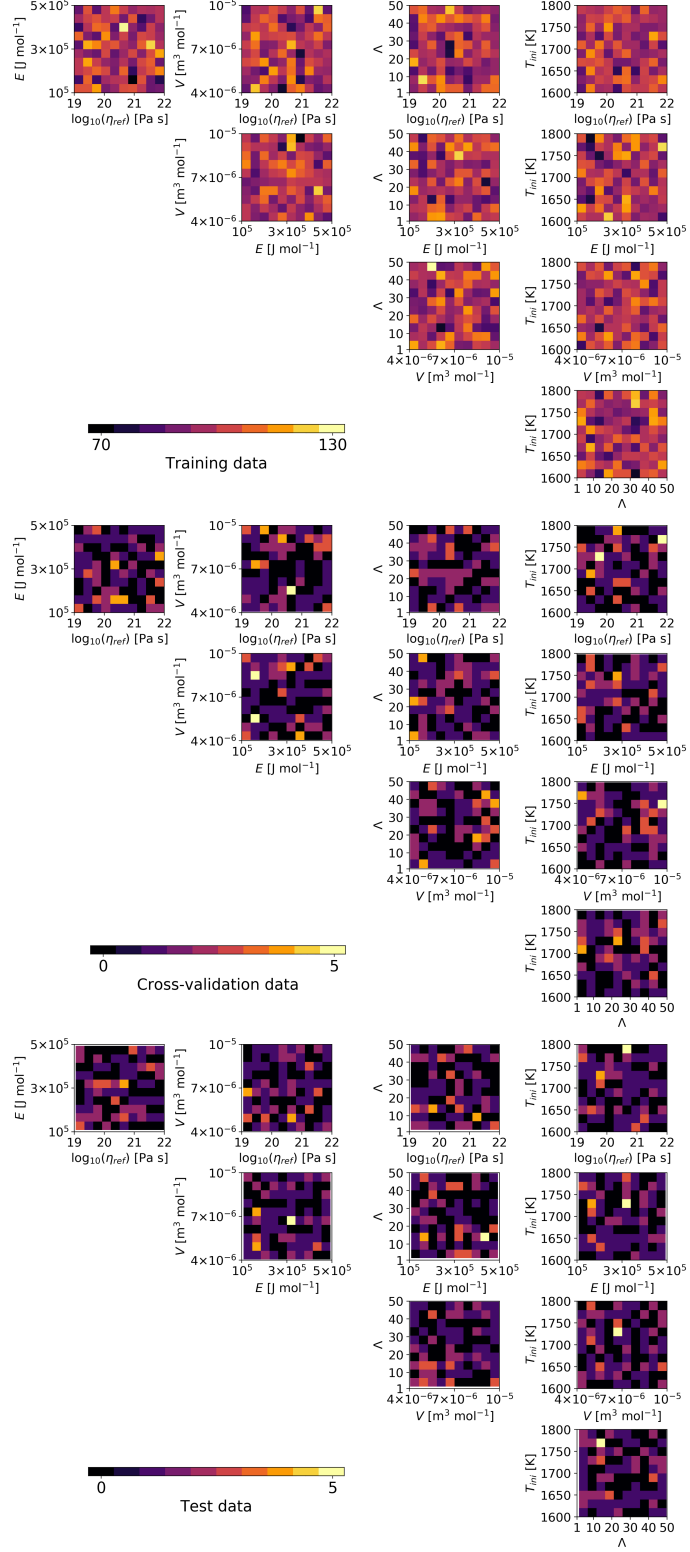


FIG. 16. Distribution of the simulation parameters in the training, cross-validation and test sets.

TABLE III. Values of fixed parameters shared by all simulations (part 2).

Parameter	Physical meaning	Value	Unit
a_0	³ Coefficient of thermal expansivity	3.15×10^{-5}	K^{-1}
a_1	³ Coefficient of thermal expansivity	1.02×10^{-8}	K^{-2}
a_2	³ Coefficient of thermal expansivity	-0.76	K
a_3	³ Coefficient of thermal expansivity	3.63×10^{-2}	GPa^{-1}
c_0	³ Coefficient of thermal conductivity	2.47	$\text{Wm}^{-1} \text{K}^{-1}$
c_1	³ Coefficient of thermal conductivity	0.33	$\text{Wm}^{-1} \text{K}^{-1} \text{GPa}^{-1}$
c_2	³ Coefficient of thermal conductivity	0.48	
e_0	⁴ Coefficient for solidus parameterization	1400	K
e_1	⁴ Coefficient for solidus parameterization	149.5	K Pa^{-1}
e_2	⁴ Coefficient for solidus parameterization	-9.4	K Pa^{-2}
e_3	⁴ Coefficient for solidus parameterization	0.313	K Pa^{-3}
e_4	⁴ Coefficient for solidus parameterization	-0.0039	K Pa^{-4}
f_0	⁵ Coefficient for liquidus parameterization	1977	K
f_1	⁵ Coefficient for liquidus parameterization	64.1	K Pa^{-1}
f_2	⁵ Coefficient for liquidus parameterization	-3.92	K Pa^{-2}
f_3	⁵ Coefficient for liquidus parameterization	0.141	K Pa^{-3}
f_4	⁵ Coefficient for liquidus parameterization	-0.0015	K Pa^{-4}

³[45] ⁴[78] ⁵[79].

where, σ is the sigmoid activation ($\sigma(x) = 1/(1 + e^{-x})$), $W_f \in \mathbb{R}^{n \times m}$ is a matrix of trainable parameters, n is the number of LSTM cells, m is the size of the input vector x_t , $U_f \in \mathbb{R}^{n \times n}$ is another matrix of trainable parameters and $b_f \in \mathbb{R}^n$ is a set of biases. The subscript f in W_f and b_f stands for the forget gate.

Then, in the ‘‘update’’ bloc, a sigmoid layer decides which values should be updated:

$$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i), \quad (\text{D2})$$

while the *SELU* layer creates new values \tilde{C}_t to be added to the state:

$$\tilde{C}_t = \text{SELU}(W_c x_t + U_c h_{t-1} + b_c). \quad (\text{D3})$$

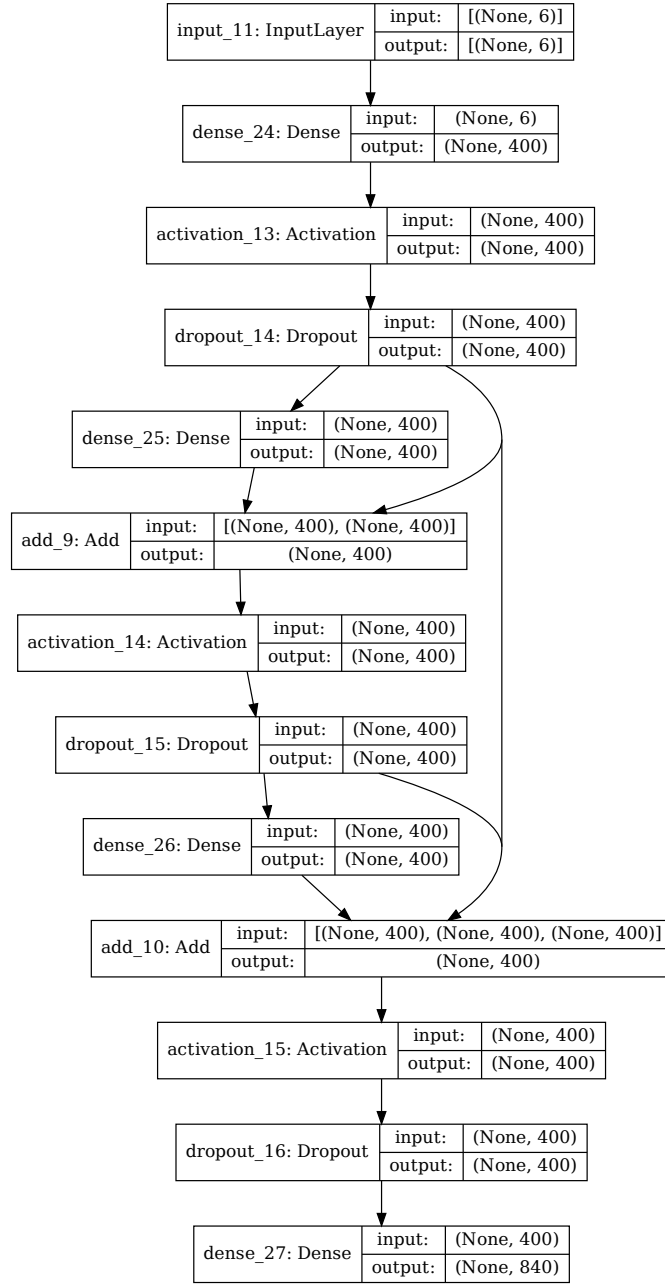


FIG. 17. A schematic of the FNN architecture with 3 hidden layers with 400 neurons each. For ease of visualization, we show only 3 hidden layers even though the trained FNN has 8.

W_i and W_c are the weights for the input connections, where subscript i denotes the weights used to update values and subscript c denoted the weights used to create new values. Similarly, U_i and U_c are weights for the recurrent connections and b_c and b_i are biases.

Using Eq. (D1)–(D3), we can now update the cell state C_t :

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t, \quad (\text{D4})$$

where, \odot is an element-wise (Hadamard) product.

Finally, a last sigmoid layer decides the amount of cell state to be outputted via the dot product of output o_t with the $SELU()$ of the cell state:

$$h_t = \sigma(W_o x_t + U_o h_{t-1} + b_o) \odot SELU(C_t), \quad (\text{D5})$$

where, W_o , U_o and b_o are the final set of trainable input weights, recurrent weights and output biases, respectively.

-
- [1] D. Breuer and W. Moore, 10.08 - dynamics and thermal history of the terrestrial planets, the moon, and io, in *Treatise on Geophysics (Second Edition)*, edited by G. Schubert (Elsevier, Oxford, 2015) second edition ed., pp. 255–305.
 - [2] P. J. Tackley, Modelling compressible mantle convection with large viscosity contrasts in a three-dimensional spherical shell using the yin-yang grid, *Physics of the Earth and Planetary Interiors* **171**, 7 (2008).
 - [3] S. Zhong, A. McNamara, E. Tan, L. Moresi, and M. Gurnis, A benchmark study on mantle convection in a 3-d spherical shell using citcoms, *Geochemistry, Geophysics, Geosystems* **9** (2008).
 - [4] M. Kronbichler, T. Heister, and W. Bangerth, High accuracy mantle convection simulation through modern numerical methods, *Geophysical Journal International* **191**, 12 (2012).
 - [5] C. Hüttig, N. Tosi, and W. Moore, An improved formulation of the incompressible navier-stokes equations with variable viscosity, *Physics of the Earth and Planetary Interiors* **220**, 11 (2013).
 - [6] N. Tosi and S. Padovan, Mercury, Moon, Mars: Surface expressions of mantle convection and interior evolution on stagnant-lid bodies, in *Mantle Convection and Surface Expressions*, edited by H. Marquardt, M. Ballmer, S. Cottar, and J. Konter (AGU Monograph Series, 2021) Chap. 17.
 - [7] M. Sambridge, Geophysical inversion with a neighbourhood algorithm-I. Searching a parameter space, *Geophysical Journal International* **138**, 479 (1999a).
 - [8] M. Sambridge, Geophysical inversion with a neighbourhood algorithm-II. Appraising the ensemble, *Geophysical Journal International* **138**, 727 (1999b).

- [9] U. Meier, A. Curtis, and J. Trampert, Global crustal thickness from neural network inversion of surface wave data, *Geophysical Journal International* **169**, 706 (2007).
- [10] P. Käüfl, A. P. Valentine, R. W. de Wit, and J. Trampert, Solving probabilistic inverse problems rapidly with prior samples, *Geophysical Journal International* **205**, 1710 (2016).
- [11] R. W. L. de Wit, A. P. Valentine, and J. Trampert, Bayesian inference of Earth's radial seismic structure from body-wave traveltimes using neural networks, *Geophys. Journal International* **195**, 408 (2013), <https://academic.oup.com/gji/article-pdf/195/1/408/1633427/ggt220.pdf>.
- [12] S. Atkins, A. P. Valentine, P. J. Tackley, and J. Trampert, Using pattern recognition to infer parameters governing mantle convection, *Physics of the Earth and Planetary Interiors* **257**, 171 (2016).
- [13] S. Agarwal, N. Tosi, P. Kessel, S. Padovan, D. Breuer, and G. Montavon, Towards constraining mars' thermal evolution using machine learning, *Earth and Space Science* **n/a**, e2020EA001484 (2021).
- [14] C. Reese, V. Solomatov, and L.-N. Moresi, Heat transport efficiency for stagnant lid convection with dislocation viscosity: Application to Mars and Venus, *Journal of Geophysical Research - Planets* **103**, 13643 (1998).
- [15] C. Dumoulin, M.-P. Doin, and L. Fleitout, Heat transport in stagnant lid convection with temperature- and pressure-dependent Newtonian or non-Newtonian rheology, *Journal of Geophysical Research* **104**, 12,759 (1999).
- [16] V. S. Solomatov and L.-N. Moresi, Scaling of time-dependent stagnant lid convection: Application to small-scale convection on earth and other terrestrial planets, *Journal of Geophysical Research* **105**, 21795 (2000).
- [17] F. Deschamps and C. Sotin, Thermal convection in the outer shell of large icy satellites, *Journal of Geophysical Research - Planets* **106**, 5107 (2001).
- [18] D. J. Stevenson, T. Spohn, and G. Schubert, Magnetism and thermal evolution of the terrestrial planets, *Icarus* **54**, 466 (1983).
- [19] M. Gurnis, A reassessment of the heat transport by variable viscosity convection with plates and lids, *Geophys. Res. Lett.* **16**, 179 (1989).
- [20] G. Schubert and T. Spohn, Thermal history of mars and the sulfur content of its core, *Journal of Geophysical Research: Solid Earth* **95**, 14095 (1990).
- [21] S. A. Hauck II, A. J. Dombard, R. J. Phillips, and S. C. Solomon, Internal and tectonic evolution of mercury, *Earth and Planetary Science Letters* **222**, 713 (2004).
- [22] J. Korenaga, Thermal evolution with a hydrating mantle and the initiation of plate tectonics in the early

- earth, *Journal of Geophysical Research: Solid Earth* **116**, 10.1029/2011JB008410 (2011).
- [23] A. Morschhauser, M. Grott, and D. Breuer, Crustal recycling, mantle dehydration, and the thermal evolution of Mars, *Icarus* **212**, 541 (2011).
- [24] N. Tosi, M. Grott, A. C. Plesa, and D. Breuer, Thermochemical evolution of Mercury's interior, *Journal of Geophysical Research: Planets* **118**, 2474 (2013b).
- [25] J. G. O'Rourke and J. Korenaga, Thermal evolution of venus with argon degassing, *Icarus* **260**, 128 (2015).
- [26] M. H. Shahnas and R. N. Pysklywec, Toward a unified model for the thermal state of the planetary mantle: Estimations from mean field deep learning, *Earth and Space Science* **7**, <https://doi.org/10.1029/2019EA000881> (2020).
- [27] S. Agarwal, N. Tosi, D. Breuer, S. Padovan, P. Kessel, and G. Montavon, A machine-learning-based surrogate model of Mars' thermal evolution, *Geophysical Journal International* 10.1093/gji/ggaa234 (2020).
- [28] T. Gillooly, N. Coltice, and C. Wolf, An anticipation experiment for plate tectonics, *Tectonics* **38**, 3916 (2019).
- [29] A.-C. Plesa, N. Tosi, M. Grott, and D. Breuer, Thermal evolution and ury ratio of mars, *Journal of Geophysical Research: Planets* **120**, 995 (2015).
- [30] A.-C. Plesa, S. Padovan, N. Tosi, D. Breuer, M. Grott, M. A. Wieczorek, T. Spohn, S. E. Smrekar, and W. B. Banerdt, The thermal state and interior structure of mars, *Geophysical Research Letters* **45**, 12,198 (2018), <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2018GL080728>.
- [31] G. Morra, D. A. Yuen, H. M. Tufo, and M. G. Knepley, Fresh Outlook in Numerical Methods for Geodynamics – Part 2: Big Data, HPC, Education, in *Encyclopedia of Geology*, edited by D. Alderton and S. A. Elias (Academic Press, 2020) pp. 841–855.
- [32] S. Pandey and J. Schumacher, Reservoir computing model of two-dimensional turbulent convection, *Phys. Rev. Fluids* **5**, 113506 (2020).
- [33] M. Raissi, P. Perdikaris, and G. E. Karniadakis, Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations, *J. Comp. Phys.* **378**, 686 (2019).
- [34] A. T. Mohan, D. Tretiak, M. Chertkov, and D. Livescu, Spatio-temporal deep learning models of 3d turbulence with physics informed diagnostics, *Journal of Turbulence* **21**, 484 (2020).
- [35] J. Masci, U. Meier, D. Cireşan, and J. Schmidhuber, Stacked convolutional auto-encoders for hierar-

- chical feature extraction, in *Artificial Neural Networks and Machine Learning – ICANN 2011*, edited by T. Honkela, W. Duch, M. Girolami, and S. Kaski (Springer Berlin Heidelberg, Berlin, Heidelberg, 2011) pp. 52–59.
- [36] X. Shi, Z. Chen, H. Wang, D.-Y. Yeung, W. kin Wong, and W. chun Woo, Convolutional lstm network: A machine learning approach for precipitation nowcasting (2015), arXiv:1506.04214 [cs.CV].
- [37] S. Bhatnagar, Y. Afshar, S. Pan, K. Duraisamy, and S. Kaushik, Prediction of aerodynamic flow fields using convolutional neural networks, *Computational Mechanics* **64**, 525–545 (2019).
- [38] S. L. Brunton, B. R. Noack, and P. Koumoutsakos, Machine learning for fluid mechanics, *Annual Review of Fluid Mechanics* **52**, 477 (2020).
- [39] G. Hirth and D. Kohlstedt, Rheology of the upper mantle and the mantle wedge: A view from the experimentalists, *AGU Monograph Series* **138**, 83 (2003).
- [40] S. D. King, C. Lee, P. E. van Keken, W. Leng, S. Zhong, E. Tan, N. Tosi, and M. C. Kameyama, A community benchmark for 2-D Cartesian compressible convection in the Earth’s mantle, *Geophysical Journal International* **180**, 73 (2010).
- [41] F. Nimmo and K. Tanaka, Early crustal evolution of mars, *Annual Review of Earth and Planetary Sciences* **33**, 133 (2005).
- [42] H. Wänke and G. Dreibus, Chemistry and accretion history of mars, *Philosophical Transactions of the Royal Society of London. Series A: Physical and Engineering Sciences* **349**, 285 (1994).
- [43] S. Padovan, N. Tosi, A.-C. Plesa, and T. Ruedas, Impact-induced changes in source depth and volume of magmatism on mercury and their observational signatures, *Nature Communications* **8** (2017).
- [44] U. Christensen and D. A. Yuen, Layered convection induced by phase transitions, *Journal of Geophysical Research: Solid Earth* **90**, 10291 (1985).
- [45] N. Tosi, D. A. Yuen, N. de Koker, and R. M. Wentzcovitch, Mantle dynamics with pressure- and temperature-dependent thermal expansivity and conductivity, *Physics of the Earth and Planetary Interiors* **217**, 48–58 (2013a).
- [46] D. Tozer, Towards a theory of thermal convection in the mantle, in *The Earth’s mantle*, edited by T. Gaskell (Academic Press, New York, 1967) pp. 327–353.
- [47] D. Breuer and T. Spohn, Early plate tectonics versus single-plate tectonics on mars: Evidence from magnetic field history and crust evolution, *Journal of Geophysical Research: Planets* **108**, <https://doi.org/10.1029/2002JE001999> (2003).
- [48] N. Tosi, A.-C. Plesa, and D. Breuer, Overturn and evolution of a crystallized magma ocean: A numerical

- parameter study for mars, *Journal of Geophysical Research: Planets* **118**, 1512 (2013).
- [49] A.-C. Plesa, N. Tosi, and D. Breuer, Can a fractionally crystallized magma ocean explain the thermochemical evolution of mars?, *Earth and Planetary Science Letters* **403**, 225 (2014).
- [50] A. Scheinberg, L. T. Elkins-Tanton, and S. J. Zhong, Timescale and morphology of martian mantle overturn immediately following magma ocean solidification, *Journal of Geophysical Research: Planets* **119**, 454 (2014).
- [51] H. H. Kieffer, Thermal model for analysis of mars infrared mapping, *Journal of Geophysical Research: Planets* **118**, 451 (2013).
- [52] A.-C. Plesa, M. Grott, N. Tosi, D. Breuer, T. Spohn, and M. A. Wieczorek, How large are present-day heat flux variations across the surface of mars?, *Journal of Geophysical Research: Planets* **121**, 2386 (2016).
- [53] J. M. Lumley, The structure of inhomogeneous turbulent flows, *Atmospheric Turbulence and Radio Wave Propagation* , 166–176 (1967).
- [54] O. Friderikos, E. Baranger, M. Olive, and D. Néron, On the stability of pod basis interpolation via grassmann manifolds for parametric model order reduction in hyperelasticity (2020), arXiv:2012.08851 [math.DG].
- [55] F. Chollet *et al.*, Keras (2015).
- [56] D. P. Kingma and J. Ba, Adam: A Method for Stochastic Optimization, arXiv e-prints , arXiv:1412.6980 (2014), arXiv:1412.6980 [cs.LG].
- [57] Joblib Development Team, Joblib: running python functions as pipeline jobs (2020).
- [58] G. Klambauer, T. Unterthiner, A. Mayr, and S. Hochreiter, Self-normalizing neural networks, *CoRR* **abs/1706.02515** (2017), arXiv:1706.02515.
- [59] S. Agarwal, N. Tosi, P. Kessel, D. Breuer, and G. Montavon, Supplementary files, <https://URLwillbeinsertedbypublisher.com> (2021).
- [60] S. Hochreiter and J. Schmidhuber, Long short-term memory, *Neural computation* **9**, 1735 (1997).
- [61] H. Eivazi, H. Veisi, M. H. Naderi, and V. Esfahanian, Deep neural networks for nonlinear model order reduction of unsteady flows, *Physics of Fluids* **32**, 105104 (2020).
- [62] C. Olah, Understanding lstm networks (2015).
- [63] Z. Che, S. Purushotham, K. Cho, D. Sontag, and Y. Liu, Recurrent neural networks for multivariate time series with missing values, *Scientific Reports* **8**, 10.1038/s41598-018-24271-9 (2018).
- [64] M. Phankokkrud and S. Wacharawichanant, A comparison of efficiency improvement for long short-

- term memory model using convolutional operations and convolutional neural network, in *2019 International Conference on Information and Communications Technology (ICOIACT)* (2019) pp. 608–613.
- [65] J. R. Fienup, Invariant error metrics for image reconstruction, *Appl. Opt.* **36**, 8352 (1997).
- [66] S. L. Brunton and J. N. Kutz, 7 data-driven methods for reduced-order modeling, in *Snapshot-Based Methods and Algorithms*, edited by P. Benner, S. Grivet-Talocia, A. Quarteroni, G. Rozza, W. Schilders, and L. M. Silveira (De Gruyter, 2020) pp. 307–344.
- [67] H. Amit, G. Choblet, P. Olson, J. Monteux, F. Deschamps, B. Langlais, and G. Tobie, Towards more realistic core-mantle boundary heat flux patterns: a source of diversity in planetary dynamos, *Progress in Earth and Planetary Science* **2** (2015).
- [68] A. C. Plesa and D. Breuer, Partial melting in one-plate planets: Implications for thermo-chemical and atmospheric evolution, *Planetary and Space Science* **98**, 50 (2014).
- [69] L. Lu, R. Pestourie, W. Yao, Z. Wang, F. Verdugo, and S. G. Johnson, Physics-informed neural networks with hard constraints for inverse design (2021), arXiv:2102.04626 [physics.comp-ph].
- [70] A. T. Mohan, N. Lubbers, D. Livescu, and M. Chertkov, Embedding hard physical constraints in neural network coarse-graining of 3d turbulence (2020), arXiv:2002.00021 [physics.comp-ph].
- [71] H. Gao, L. Sun, and J.-X. Wang, Phygeonet: Physics-informed geometry-adaptive convolutional neural networks for solving parameterized steady-state pdes on irregular domain, *Journal of Computational Physics* **428**, 110079 (2021).
- [72] F. Brockherde, L. Li, K. Burke, and K.-R. Müller, By-passing the kohn-sham equations with machine learning, *Nature Communications* **8** (2017).
- [73] B. Hamzi, R. Maulik, and H. Owhadi, Data-driven geophysical forecasting: Simple, low-cost, and accurate baselines with kernel methods (2021).
- [74] N. Margenberg, C. Lessig, and T. Richter, Structure preservation for the deep neural network multigrid solver (2020), arXiv:2012.05290 [math.NA].
- [75] S. Agarwal, Supporting code for deep learning for surrogate modelling of two-dimensional mantle convection (2021).
- [76] A. Brand, L. Allen, M. Altman, M. Hlava, and J. Scott, Beyond authorship: Attribution, contribution, collaboration, and credit, *Learned Publishing* **28** (2015).
- [77] G. Schubert, D. L. Turcotte, and P. Olson, *Mantle convection in the Earth and planets* (Cambridge University Press, 2001).
- [78] C. Herzberg, P. Raterron, and J. Zhang, New experimental observations on the an-

- hydrous solidus for peridotite klb-1, *Geochemistry, Geophysics, Geosystems* **1** (2000), <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2000GC000089>.
- [79] J. Zhang and C. Herzberg, Melting experiments on anhydrous peridotite klb-1 from 5.0 to 22.5 gpa, *Journal of Geophysical Research - Solid Earth* **99**, 17729 (1994), <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/94JB01406>.
- [80] P. Van Keken, Cylindrical scaling for dynamical cooling models of the earth, *Physics of the Earth and Planetary Interiors* **124**, 119 (2001).
- [81] H. Wänke, G. Dreibus, S. K. Runcorn, G. Turner, and M. M. Woolfson, Chemical composition and accretion history of terrestrial planets, *Philosophical Transactions of the Royal Society of London. Series A, Mathematical and Physical Sciences* **325**, 545 (1988), <https://royalsocietypublishing.org/doi/pdf/10.1098/rsta.1988.0067>.