



Data-driven physics-constrained recurrent neural networks for multiscale damage modeling of metallic alloys with process-induced porosity

Shiguang Deng^{1,2} · Shirin Hosseinmardi³ · Libo Wang¹ · Diran Apelian¹ · Ramin Bostanabad³

Received: 3 May 2023 / Accepted: 2 December 2023

© The Author(s), under exclusive licence to Springer-Verlag GmbH Germany, part of Springer Nature 2024

Abstract

Computational modeling of heterogeneous materials is increasingly relying on multiscale simulations which typically leverage the homogenization theory for scale coupling. Such simulations are prohibitively expensive and memory-intensive especially when modeling damage and fracture in large 3D components such as cast metallic alloys. To address these challenges, we develop a physics-constrained deep learning model that surrogates the microscale simulations. We build this model within a mechanistic data-driven framework such that it accurately predicts the effective microstructural responses under irreversible elasto-plastic hardening and softening deformations. To achieve high accuracy while reducing the reliance on labeled data, we design the architecture of our deep learning model based on damage mechanics and introduce a new loss component that increases the thermodynamical consistency of the model. We use mechanistic reduced-order models to generate the training data of the deep learning model and demonstrate that, in addition to achieving high accuracy on unseen deformation paths that include severe softening, our model can be embedded in 3D multiscale simulations with fracture. With this embedding, we also demonstrate that state-of-the-art techniques such as teacher forcing result in deep learning models that cause divergence in multiscale simulations. Our numerical experiments indicate that our model is more accurate than pure data-driven models and is much more efficient than mechanistic reduced-order models.

Keywords Multiscale damage modeling · Recurrent neural networks · Data-driven surrogate · Path dependency · Physics constraints

1 Introduction

Heterogeneous materials are increasingly used in many engineering applications. Analyzing the behavior of such materials often relies on multiscale simulations such as the FE² method [1] which is a popular homogenization-based concurrent multiscale model that uses the finite element method (FEM) at two spatial scales. Despite the recent advancements in software/hardware and mechanics theory [2], the simulation of hierarchical materials via FE² is still prohibitively

costly. Consider the multiscale model in Fig. 1a where each integration point (IP) of the macroscale component represents a microstructure with complex local morphologies. In this model, the two-scale spatial discretization requires large memory storage and also results in long runtimes since the solver repeatedly iterates between the scales. These challenges are exacerbated in the presence of microstructural deformations that are path-dependent and involve damage. That is, the evaluation of the microstructural responses is the primary computational bottleneck. Our goal in this paper is to address such bottlenecks by developing a deep learning (DL) model that surrogates the microstructural analyses in 3D multiscale simulations that involve plasticity and fracture.

The computational costs of microstructure analyses is especially high in the presence of fracture whose modeling is one of the most important sub-branches of solid mechanics. Fracture mechanics aims to quantitatively study the initiation, accumulation, and propagation of damage in materials and its numerical simulation is generally carried out by two

✉ Ramin Bostanabad
Raminb@uci.edu

¹ Materials Science and Engineering, University of California, Irvine, USA

² Department of Mechanical Engineering, Northwestern University, Evanston, USA

³ Department of Mechanical and Aerospace Engineering, University of California, Irvine, USA

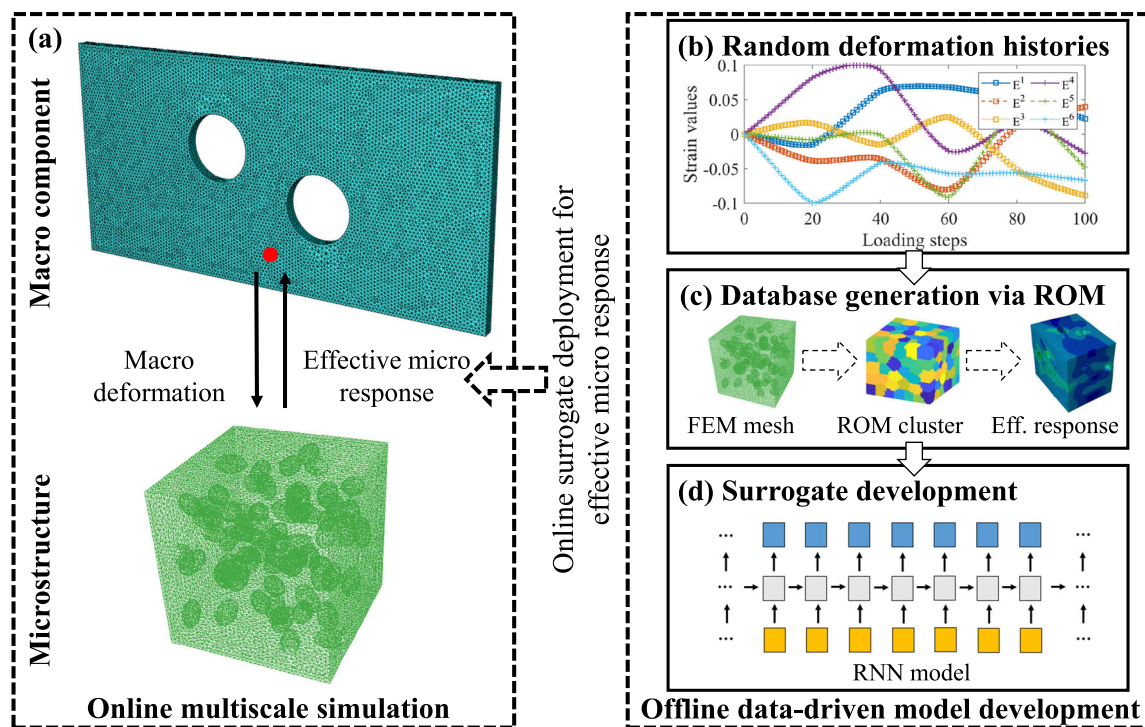


Fig. 1 Multiscale simulations: **a** Each integration point of the macroscale component represents a microstructure whose effective response is needed as the solver iterates between the two scales. We build a physics-constrained data-driven material model to surrogate the expensive microscale analyses amid online multiscale simulations. Our

data-driven framework has three major offline components: **b** deformation representation, **c** response database generation via mechanistic reduced-order models (ROMs), and **d** physics-constrained deep learning

approaches [3]: a discrete approach and a continuous one. The former approach explicitly models discontinuous displacement fields across fracture interfaces. Some of the early discrete models include linear elastic fracture mechanics (LEFM) [4] and cohesive zone models (CZM) [5] which require additional criteria to formulate the crack initiation, growth, propagation, and branching. To address the strong discontinuity-induced numerical difficulty in a standard FEM with continuous displacement fields, recent works consider LEFM with nodal duplication and remeshing techniques [6], CZM with zero-thickness cohesive surfaces [7], enriched FEM with elemental and nodal enrichment [8], and the extended FEM (XFEM) [9, 10] where the approximation functions are enriched and additional degrees of freedom are provided to represent the discontinuity without mesh refinement.

The continuous approach models fracture as a stress/modulus degradation process via strain softening on a continuous displacement field. The two most popular continuous models are the smeared crack model [11, 12] (which smears the displacement discontinuity over a fracture zone of finite width) and the continuum damage mechanics [13, 14] (which considers the influence of micro defects on damage evo-

lution based on phenomenological damage parameters). Continuum damage mechanics suffers from severe mesh dependency due to the lack of objectivity of softening constitutive formulation which causes erroneous imaginary wave speed amid damage propagation. Remedies include viscous regularization [15], cosserat micropolar theory [16], nonlocal continuum theory [17, 18], and gradient-enhanced damage models [19, 20].

As alternatives to expensive methods such as the FE^2 and direct numerical simulations (DNS), mechanistic reduced-order models (ROMs) are developed which can significantly accelerate computational plasticity and damage mechanics. The main idea behind ROMs is to reduce the number of unknown variables (e.g., stresses, strains, or internal variables such as the damage parameters) while striking a balance between accuracy and efficiency. For example, the transformed field analysis method [21] and its non-uniform variant [22] employ proper orthogonal decomposition to reduce material state variables by expressing arbitrary strain fields as a subspace representation of pre-computed eigenstrains. Clustering-based ROMs reduce unknown variables by agglomerating a large number of IPs into a few clusters. For instance, the self-consistent clustering [23] and its variant

the virtual clustering analysis [24] methods assume IPs with similar elastic responses behave similarly during inelastic deformations and solve incremental Lippmann-Schwinger equations to approximate the evolution of cluster-wise material responses. Deflated clustering analysis (DCA) [25] utilizes clusters to decompose both macroscale and microscale domains where macro analysis is faithfully accelerated in a deflation space while the effective microstructural responses are approximated in a coarse-grained manner where close-by IPs are presumed to share the same behaviors. DCA's robustness and efficiency are further improved in [26] where both spatial and temporal dimensions are adaptively reduced for elasto-plastic deformations with softening. While ROMs dramatically accelerate multiscale simulations, their runtimes are still quite high (especially in the presence of softening). Additionally, ROMs lack solution transferability in that the expensive data of one instance of the model is not reused (e.g., the full strain–stress history obtained for a microstructure corresponding to a macroscopic IP is not reused in another multiscale simulation).

Machine learning (ML) provides a feasible avenue for building transferable and extremely fast material models. Among the various ML techniques such as Gaussian processes (GPs) [27–30] and polynomial chaos expansion [31], neural networks (NNs) are increasingly employed in computational solid mechanics to build data-driven material models [32, 33]. For instance, NNs are used to learn visco-plasticity [34], cyclic plasticity [35], interface mechanism [36], and anisotropic electrical behaviors [37]. More recently, Mianroodi et al. [38] build an NN to calculate local stress distributions in non-homogeneous microstructures with elasto-plastic behaviors. Haghight et al. [39] incorporate the momentum balance and constitutive relations into a feed-forward NN and demonstrate the improved extrapolation capability for single-scale elasto-plastic simulations. Peivaste et al. [40] develop a convolutional NN to surrogate costly phase-field simulations to learn microstructural grain evolution.

Data-driven material models are increasingly built via recurrent neural networks (RNNs) to learn path-dependent constitutive laws that govern elasto-plastic deformations. For example, Mozaffar et al. [41] successfully use an RNN to learn plasticity with distortional hardening on 2D fiber composite microstructures. Wang et al. [42] develop an RNN to link information from different scales via recursive homogenization to capture the multiscale hydro-mechanical coupling effects of heterogeneous media with various pore sizes. In these works, various methods are proposed for data generation. For instance, Wu et al. [43] design a random walk algorithm to generate a database of effective elasto-plastic hardening behaviors under cyclic and non-proportional loading paths. An on-demand sampling strategy is adopted by Ghavamian et al. [44] that reduces sampling space by run-

ning prior macro models to collect the strain–stress sequences for the subsequent RNN's learning process. This strategy reduces sampling efforts and improves prediction accuracy but reduces the generalization power since the trained model can be only applied to the macro component that is used to collect the training sequences. In a recent work [45], Logarzo et al. use an RNN to learn the hardening behavior of a 2D composite microstructure under a wide range of deformation histories that are sampled from the space of principal strains.

All aforementioned RNN surrogates are black-box or pure data-driven models whose accuracy relies on large training datasets. Building such datasets is very challenging for 3D microstructural analyses that involve softening. Since infusing physical laws into the training process can improve the reliance on data and energy consistency, we rigorously explore this direction to derive the physical constraints that must be met when surrogating elasto-plastic deformations with softening. That is, our main contribution is to develop a physics-constrained RNN that surrogates the micro analyses amid online multiscale simulations that involve softening. Compared to the reviewed ROMs and pure data-driven models, our surrogate is computationally efficient, memory-light, physics consistent, and transferable.

The rest of the paper is organized as follows. In Sect. 2, we briefly review the homogenization-based concurrent multiscale modeling. In Sect. 3, we propose our physics-constrained data-driven model to surrogate effective elasto-plastic microstructural responses that may involve damage. Specifically, we derive two constraints based on the continuum damage mechanics and integrate them in our surrogate to reduce data reliance and improve prediction accuracy. In Sect. 4, we illustrate the efficiency and accuracy of our data-driven model by (1) evaluating its predictions of microstructural effective responses subject to random deformation paths, and (2) embedding it in a number of multiscale structures subject to complex cyclic loading conditions with hardening and softening material behaviors. We conclude our paper with some notes on the contributions, limitations, and future research directions in Sect. 5.

2 Review of Homogenization-Based Multiscale Modeling

Our multiscale damage analysis is based on the first-order homogenization model which we review in this section. We also briefly review the continuum damage modeling in Appendix A, and in Appendix B we discuss a constitutive hybrid integration scheme that we previously developed to address softening-induced numerical instability.

The first-order computational homogenization assumes scale separation between a macroscale component and its microscopic features. In solving multiscale systems, the solu-

tions at the macroscale and microscale are coupled via the Hill-Mandel condition [46]. This condition equates the density of virtual internal work of a macroscale IP to the volume average of the virtual work in the associated microstructure subject to kinematically admissible displacement fields:

$$\mathbf{S}_M : \delta \mathbf{E}_M = \frac{1}{|\Omega_{0m}|} \int_{\Omega_{0m}} \mathbf{S}_m : \delta \mathbf{E}_m d\Omega \quad (1)$$

where \mathbf{S}_M , $\delta \mathbf{E}_M$, \mathbf{S}_m and $\delta \mathbf{E}_m$ represent the macroscopic stress, virtual macroscopic strain, microscopic stress, and virtual microscopic strain, respectively. The subscripts M and m indicate the macroscale and microscale, respectively. The $:$ operator represents the double dot product contracting a pair of repeated indices. In addition, Ω_{0m} and $|\Omega_{0m}|$ indicate the reference microstructural domain and its volume, respectively. Following the virtual energy condition in Equation (1), the macroscopic effective stress and virtual strain can be expressed as the volume average of their micro counterparts as:

$$\mathbf{S}_M = \frac{1}{|\Omega_{0m}|} \int_{\Omega_{0m}} \mathbf{S}_m d\Omega; \quad \delta \mathbf{E}_M = \frac{1}{|\Omega_{0m}|} \int_{\Omega_{0m}} \delta \mathbf{E}_m d\Omega \quad (2)$$

The stress and strain at both the macroscale and microscale need to satisfy equilibrium equations at the corresponding length scale. For instance, under the infinitesimal deformation assumption, the macro-solutions at the arbitrary macroscopic IP \mathbf{P} can be computed by solving the following boundary value problem (BVP):

$$\nabla_0 \cdot \mathbf{S}_M(\mathbf{P}) + \mathbf{b}_M = \mathbf{0} \quad \forall \mathbf{P} \in \Omega_{0M} \quad (3a)$$

$$\mathbf{u}_M(\mathbf{P}) = \bar{\mathbf{u}}_M \quad \forall \mathbf{P} \in \Gamma_{0M}^D \quad (3b)$$

$$\mathbf{S}_M(\mathbf{P}) \cdot \mathbf{n}_M = \bar{\mathbf{t}}_M \quad \forall \mathbf{P} \in \Gamma_{0M}^N \quad (3c)$$

where \mathbf{u}_M is the unknown macroscopic displacement in Ω_{0M} and $\bar{\mathbf{u}}_M$ is the prescribed displacement on the Dirichlet boundary Γ_{0M}^D over the undeformed macroscopic domain Ω_{0M} with an outward unit vector \mathbf{n}_M . Also, ∇_0 indicates the gradient operator with respect to the original configuration. \mathbf{b}_M and $\bar{\mathbf{t}}_M$ represent the body force and prescribed surface traction on the Neumann boundary Γ_{0M}^N , respectively.

In a similar manner, the strong form of the microscale equilibrium equations can be written as a BVP for the microstructure or representative volume element (RVE) composed of micro IPs \mathbf{p} as:

$$\nabla_0 \cdot \mathbf{S}_m(\mathbf{p}) = \mathbf{0} \quad \forall \mathbf{p} \in \Omega_{0m} \quad (4a)$$

$$\mathbf{S}_m(\mathbf{p}) \cdot \mathbf{n}_m = \bar{\mathbf{t}}_m \quad \forall \mathbf{p} \in \Gamma_{0m} \quad (4b)$$

where $\bar{\mathbf{t}}_m$ indicates the surface traction per unit area over the reference microstructural boundary Γ_{0m} with an outward unit normal vector \mathbf{n}_m .

3 Proposed Physics-Constrained Data-Driven Surrogate

To reduce the computational costs of multiscale simulations that involve hardening and softening, we follow the data-driven framework that is schematically illustrated in Fig. 1b–d. Our framework uses the following three modules to build a physics-informed material model that surrogates the nested microstructural analyses. The first two modules create the training dataset for the material model and the last module builds an RNN using the generated data and domain knowledge.

- Module 1: Exploration of the deformation space.** Deformation paths are extremely high dimensional as they have a sequential nature. To efficiently explore such a high-dimensional space, we utilize random processes and design of experiments (DoE). This exploration is an extremely important step because in a multiscale simulation each macro IP (and hence its corresponding RVE) undergoes a unique deformation path that depends on the IP's location, the material properties, the applied macro boundary and initial conditions, and the geometry of the macro component.
- Module 2: Response collection.** The first major computational bottleneck of our framework is obtaining the microstructural responses to the deformation paths generated in Module 1. These high costs are primarily associated with simulating softening and we address them by using ROMs which leverage hybrid time integration to avoid softening-induced solver divergence.
- Module 3: Physics-informed surrogate modeling.** Fitting an accurate surrogate to the generated training data is a challenging and time-consuming process since the generalization power¹ of the RNN strongly depends on its architecture, training data size, and the training mechanism. We use domain knowledge to dramatically reduce the sensitivity of the RNN to these factors such that it can be used as a transferable constitutive law in a wide range of multiscale simulations.

We describe these three modules in more detail below. In Sect. 3.1 we elaborate on the data generation process which builds a set of independent and systematically sampled microstructural deformation-response sequences. This

¹ The accuracy of the surrogate in predicting the microstructural response given deformation paths that are not seen in training.

dataset is then used in Sect. 3.2 to train an RNN that serves as the data-driven material model at the microscale. To improve this model’s accuracy on unseen deformation paths, we incorporate two physics constraints in Sect. 3.3. In Sect. 3.4, we show the integration procedure of our surrogate in multiscale solvers.

3.1 Design of Experiments

While using domain knowledge reduces the reliance on training data, building a transferable² RNN is still a data intensive and computationally expensive task since thousands of costly samples are needed. Hence, it is important to ensure that the training data provides as much information as possible by maximally exploring the deformation space. Below, we first define this space and then describe our sampling mechanism.

In the DoE, we assume every strain path starts from a relaxing state (with zero initial strains and no residual stresses) and evolves to a final state by n_{load} load steps. Additionally, we presume that (1) the maximum strain value in any direction and at any load step is smaller than the user-defined threshold ζ_1 , and (2) our material’s bulk modulus is fairly large such that the deformation-induced material volume change is within the user-defined limit ζ_2 . We use these two practical constraints to reduce the sampling space and we express them as:

$$|E_n^i| \leq \zeta_1; \quad |E_n^{vol}| \leq \zeta_2 \tag{5}$$

where E_n^i represents the i^{th} component of the strain vector at load step $n \in \{1, 2, \dots, n_{load}\}$, $i \in \{1, 2, 3, 4, 5, 6\}$ indicates the six components of 3D strains in which $i = \{1, 2, 3\}$ designate normal strains and $i = \{4, 5, 6\}$ represent shear strain components. We note that $E_n^{vol} = E_n^1 + E_n^2 + E_n^3$ is the volumetric strain standing for the material volume change after deformation.

Without loss of generality, we require all deformation paths to have n_{load} load steps. We characterize each dimension of a load path³ with n_c evenly spaced control points whose strain values are sampled via a space-filling algorithm such as the Sobol sequence. To obtain the strain values at all n_{load} load steps, we use a one-dimensional zero-mean GP to interpolate the values assigned to the control points. The correlation function of this GP is:

$$r(n, n') = \exp(-w(n - n')^2) \tag{6}$$

where n and n' are two load steps, w is the scale parameter that controls the roughness of the interpolated curve, and the

² By Transferable we mean an RNN that can be used as the constitutive law at all macro IPs in a wide range of multi-scale simulations.

³ In 3D, a load path consists of 6 strain sequences where each sequence is of length n_{load} .

exponent 2 ensures that the generated path is differentiable. In our dataset, we change the value of w across the DoE samples to increase the variability in the generated paths. That is, we use a single Sobol sequence to generate both w and the values of all strain components at all control points.

To ensure the values of all six strain components at the control points satisfy the two constraints in Eq. 5, we generate a large DoE via the Sobol sequence (which is extremely fast) and then select n_p valid points from it. Specifically, we choose the independent variables whose bounds are known as the DoE dimensions, i.e., E_c^{vol} and E_c^i with $i \in \{1, 2, 4, 5, 6\}$. Then, we find the third normal strain for the control point c via $E_c^3 = E_c^{vol} - E_c^1 - E_c^2$.

We create a total of n_p deformation paths where each one represents the temporal evolution of the six independent strain components. We plot ten example strain paths in Fig. 2 which demonstrates that the random shear strains span the entire hypercube-shaped deformation space constrained by ζ_1 while the normal strain components are additionally confined between the two hyper-planes that represent the volumetric strain constraint defined by ζ_2 . We also plot the 2D projections of the random strain sequences in Fig. 2 to show that the normal and shear components start from the relaxing state without any strain values. In addition, we observe that the highly complex deformation histories consist of multiple loading-unloading-reloading cycles.

After we generate the random strain paths, we use ROMs to compute the microstructural effective responses for each strain sequence. Specifically, we impose the microstructural displacement boundary conditions by the affine boundary condition as:

$$\mathbf{u}_m(\mathbf{p}) = \mathbf{E}_M \Delta \mathbf{p} \quad \forall \mathbf{p} \in \Gamma_{0m} \tag{7}$$

where the microstructural displacement boundary condition \mathbf{u}_m depends on the macro strain tensors \mathbf{E}_M (generated from GP interpolations) and the relative coordinates $\Delta \mathbf{p}$ of the nodes on the microstructural boundary Γ_{0m} . From this BVP, we proceed to solve the microstructural local stress \mathbf{S}_m , and compute the effective stress \mathbf{S}_M via Eq. 2.

In Sect. 4.1.1 we detail the specific values of ζ_1 and ζ_2 (user-defined DoE sampling constraints), n_c (number of control points for random strain), n_p (number of deformation paths), and n_{load} (number of sequential loading steps) that we use in our experiments.

3.2 Vanilla Data-Driven Surrogate

In this section we first review the working principles of RNNs and then explain how to use the generated data in Sect. 3.1 to train a vanilla data-driven RNN (in Sect. 4.1.2 we compare the accuracy of this RNN to the physics-constrained RNN

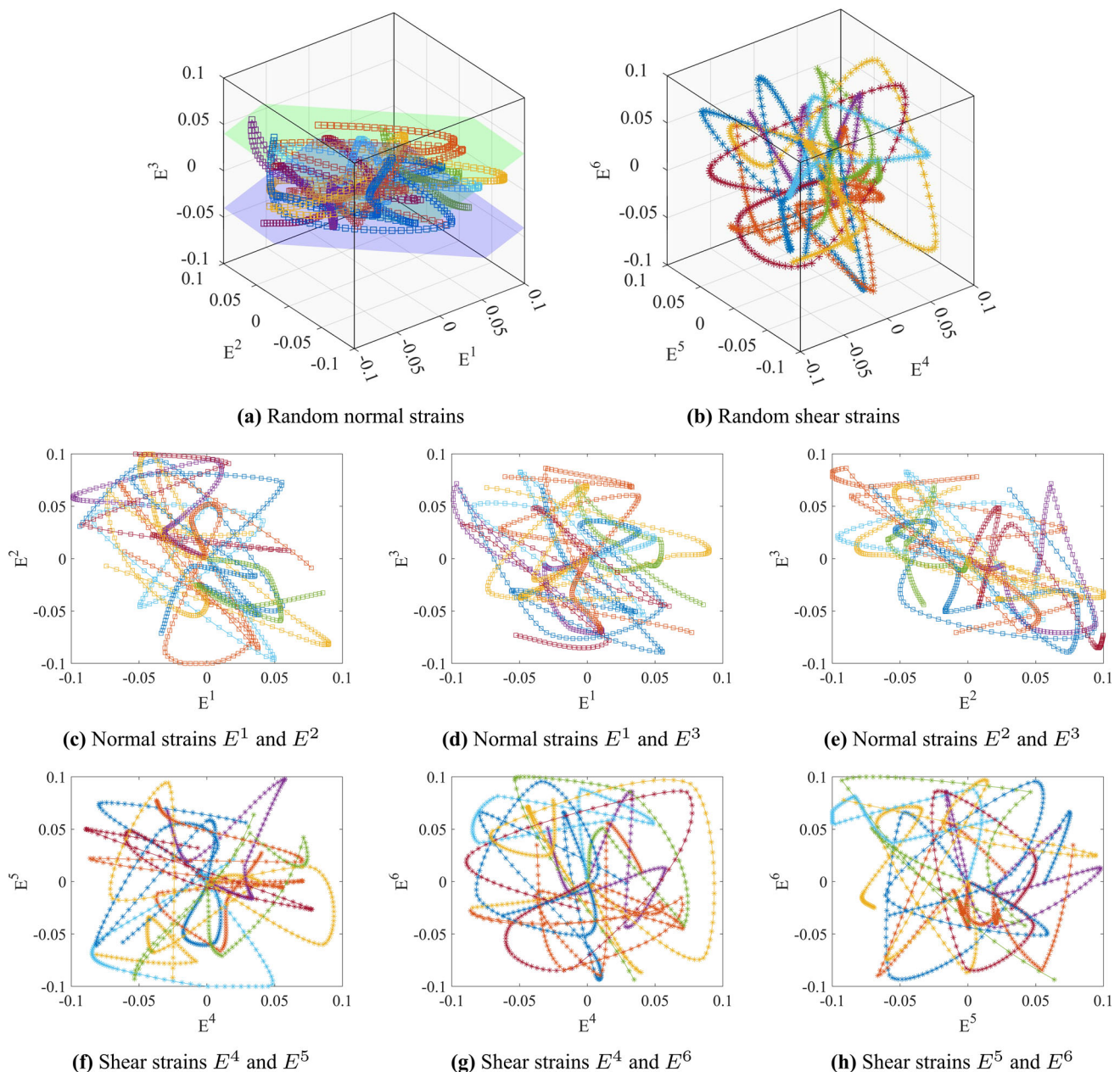


Fig. 2 Example deformation paths in 3D simulations: Ten random strain paths are illustrated. The normal and shear strain components are shown in (a) and (b), respectively. The 2D projections of these strain paths are provided in (c)–(e) for normal components and in (f)–(g) for shear components

of Sect. 3.3). We also discuss the limitations of excluding domain knowledge from the training process.

RNN is a special type of NN that is designed to learn from sequences (e.g., time series data) which cannot be accurately learned by basic NNs such as feed-forward neural networks [47, 48] (FFNNs, aka multi-layer perceptrons). FFNNs are a collection of neurons arranged in multiple layers such that each neuron has one-way connections to the neurons of the subsequent layer. In an FFNN, each neuron performs a relatively simple mathematical operation where a (typically)

nonlinear activation function is applied to the summation of a bias term and a weighted sum of the neuron's inputs as:

$$\mathbf{x}^l = f(\mathbf{W}^l \mathbf{x}^{l-1} + \mathbf{b}^l) \quad (8)$$

where f is a vector of activation functions (e.g., hyperbolic tangent, rectified linear unit or ReLU, leaky ReLU, or swish), \mathbf{x}^{l-1} are the outputs of the previous layer's neurons, \mathbf{W}^l and \mathbf{b}^l are the weight matrix and bias vector of the layer l , respectively, and \mathbf{x}^l are the outputs of layer l .

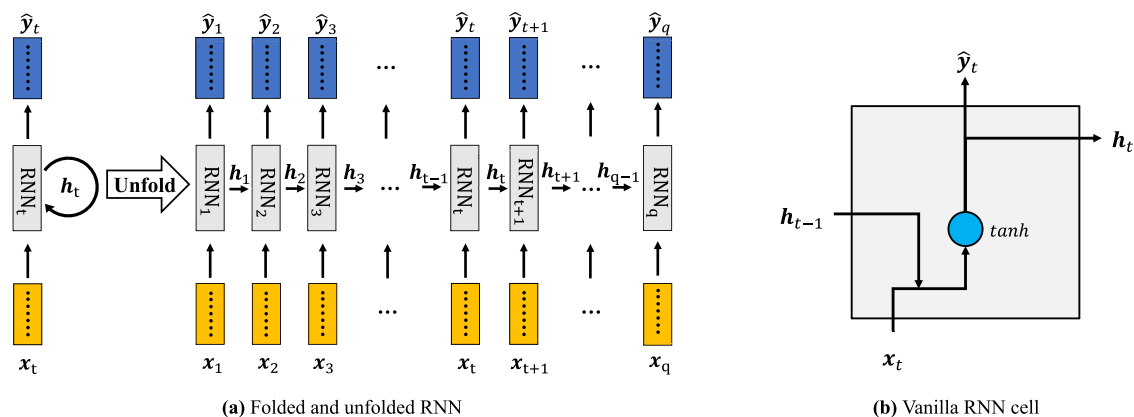


Fig. 3 Computational graph of vanilla data-driven RNN: **a** Folded and unfolded representations of an RNN that maps a sequence of inputs to a sequence of outputs. h_t is the state variable and captures the effects of the past on the present. **b** Internal structure and data flow in an RNN

Operations such as the one in Equation (8) cannot efficiently learn from sequences as their structure does not have a mechanism to leverage the past in predicting the current response. RNNs [49] address this issue via the so-called state variables⁴ that capture the effects of the past events (i.e., past inputs/outputs in the sequence) on the current response. To demonstrate this, we draw an RNN cell and its equivalent unfolded computational graph in Fig. 3a where the RNN cell relates the input sequence x_t to a series of outputs y_t with t representing the pseudo-time (equivalent to our load steps). The mathematical operations that an RNN cell performs at time step t are schematically demonstrated in Fig. 3b and read as:

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h) \tag{9a}$$

$$\hat{y}_t = W_{hy}h_t + b_y \tag{9b}$$

where the hidden state h_t depends on the current input state x_t and the previous hidden state h_{t-1} . In addition, W_{xh} , W_{hh} , and W_{hy} are weighting matrices corresponding to input-to-hidden, hidden-to-hidden, and hidden-to-output affine transformations, respectively. b_h and b_y are the bias terms associated with h_t and estimated outputs, respectively.

As indicated in Fig. 3, regardless of the sequence length, an RNN always has the same input size because it is specified in terms of the transition from one state to another state rather than being specified in terms of a variable-length history of states. That is, the same transition function with the same parameters is used at every time step. This parameter

⁴ Interestingly, state variables are also used in classical constitutive laws such as the over-stress tensor that is used when modeling multi-axial large-strain kinematic hardening.

cell at time instance t where a hyperbolic tangent function maps the weighted current inputs and hidden variables from the previous time step to the current outputs and hidden variables

sharing provides RNNs with major advantages over FFNNs in sequence learning.

RNNs are data-intensive models especially when learning complex and long sequences. They also suffer from numerical issues such as vanishing and exploding gradients [50] that prevent the RNN from learning long-range dependencies. To address these issues, more advanced cells such as long short-term memory (LSTM) [51] and gated recurrent unit (GRU) [52] are developed. In this work, we employ GRU cells and review their structure in Appendix D.

While GRU cells are more efficient than vanilla RNN cells, they still need large training data to achieve *sufficient* accuracy. This sufficiency condition is driven by our application, that is, our RNN should be accurate enough such that a multi-scale simulation converges to the ground truth when the RNN is used as the microstructural material model at all macro IPs. In addition to being data-intensive, the performance of RNNs is sensitive to factors such as the architecture and training mechanism (e.g., batch size, learning rate, regularization weight, etc.). To reduce the reliance on data and sensitivity to these factors, in Sect. 3.3 we propose to use domain knowledge in designing the architecture and loss function of the RNN.

3.3 Physics-Constrained Surrogate

The samples in our training dataset are expensive since obtaining the microstructural response under a deformation path requires running 3D elasto-plastic simulations with damage. Hence, we cannot afford to build a very large training dataset which, in turn, challenges building an accurate RNN. To address this issue, we leverage domain knowledge as model constraints to design the architecture and loss function of our RNN. Our additions allow to reduce

the number of trials that an analyst has to test before finding a good model. This reduction is particularly advantageous since training RNNs relies on back-propagation through time (BPTT) which cannot be parallelized (as the forward propagation graph is inherently sequential, i.e., each time step can only be predicted after the previous one is already computed) and is also memory intensive (since all the values computed in the forward pass must be stored until they are reused during the backward pass).

3.3.1 Soft Constraint: Thermodynamics

Let us assume an arbitrary micro point in a microstructural analysis is subject to an iso-thermal elasto-plastic deformation. We can compute its total work rate per unit volume \dot{W} via thermodynamics principles [53] as:

$$\dot{W} = \dot{\psi} + \Phi \quad (10)$$

where $\dot{\psi}$ represents the rate of Helmholtz free energy and Φ accounts for the rate of dissipated energy including the dissipation from plasticity, damage, damping, etc. For general elasto-plastic material behaviors, we can decompose the rate of work into elastic and plastic parts:

$$\dot{W} = \dot{W}^{el} + \dot{W}^{pl} \quad (11)$$

where the elastic work rate \dot{W}^{el} of the micro IP is equal to the rate of recoverable elastic free energy or strain energy $\dot{\psi}^{el}$, while the plastic work rate \dot{W}^{pl} is equal to the sum of the conditionally recoverable plastic free energy $\dot{\psi}^{pl}$ and the irrecoverable dissipation rate [54]. That is:

$$\dot{W}^{el} = \dot{\psi}^{el} = \mathbf{S}_m : \dot{\mathbf{E}}_m^{el}; \quad \dot{W}^{pl} = \dot{\psi}^{pl} + \Phi = \mathbf{S}_m : \dot{\mathbf{E}}_m^{pl} \quad (12)$$

We write the total rate of work per unit volume \dot{W} of an RVE as the multiplication of the micro stress \mathbf{S}_m and the rate of microscopic total strain $\dot{\mathbf{E}}_m$:

$$\dot{W} = \mathbf{S}_m : \dot{\mathbf{E}}_m^{el} + \mathbf{S}_m : \dot{\mathbf{E}}_m^{pl} = \mathbf{S}_m : \dot{\mathbf{E}}_m \quad (13)$$

where we use the additive decomposition rule for the strain. We compute the total work by integrating the work rate over the time interval and spatial domain Ω_{0m} . Additionally, we compute the total work by invoking the Hill-Mandel energy condition from Eq. 1 while assuming the strain rates are in the hyperspace of the virtual strains:

$$\begin{aligned} W &= \int_t \int_{\Omega} \dot{W} d\Omega_{0m} dt = \int_t \int_{\Omega} \mathbf{S}_m : \dot{\mathbf{E}}_m d\Omega_{0m} dt \\ &= |\Omega_{0m}| \int_t \mathbf{S}_M : \dot{\mathbf{E}}_M dt = |\Omega_{0m}| \int_t \mathbf{S}_M d\mathbf{E}_M \end{aligned} \quad (14)$$

We now decompose this total work to its components to find the constraints that the effective macroscale stress and strain fields should satisfy. We write W as a summation of total strain energy W^{el} and total plastic work W^{pl} :

$$W = W^{el} + W^{pl} \quad (15)$$

Assuming linear elasticity, we can show the total strain energy of the RVE as:

$$\begin{aligned} W^{el} &= \int_t \int_{\Omega} \mathbf{S}_m : \dot{\mathbf{E}}_m^{el} d\Omega_{0m} dt = \int_{\Omega} \int_{\mathbf{E}_m^{el}} \mathbf{S}_m d\mathbf{E}_m^{el} d\Omega_{0m} \\ &= \frac{1}{2} \int_{\Omega} \mathbf{E}_m^{el} : \mathbb{C}_m^{el} : \mathbf{E}_m^{el} d\Omega_{0m} \geq 0 \end{aligned} \quad (16)$$

where \mathbb{C}_m^{el} represents the elastic modulus of a micro IP which is equal to $(1 - D_m)\mathbb{C}^{el}$ if damage occurs (with a micro damage parameter D_m) and \mathbb{C}^{el} if there is no damage. Since $0 \leq D_m \leq 1$, it is straightforward to see $W^{el} \geq 0$.

Similarly, we can compute W^{pl} by spatiotemporally integrating \dot{W}^{pl} , and it is equal to the sum of total dissipated energy and total plastic free energy as:

$$W^{pl} = W^{di} + W^{pf} \quad (17)$$

where W^{di} can be expressed as the spatiotemporal integration of the non-negative dissipation rate:

$$W^{di} = \int_t \int_{\Omega} \Phi d\Omega_{0m} dt \geq 0 \quad (18)$$

where the non-negativity is due to the fact that $\Phi \geq 0$. The total plastic free energy equals the integrated rate of plastic free energy:

$$W^{pf} = \int_t \int_{\Omega} \dot{\psi}^{pl} d\Omega_{0m} dt = \int_{\Omega} \psi^{pl} d\Omega_{0m} \quad (19)$$

where ψ^{pl} stands for the density of the plastic free energy in the RVE and it can be decomposed into isotropic and anisotropic parts [54] as:

$$\psi^{pl} = \psi_{iso}^{pl} + \psi_{ani}^{pl}; \quad \psi_{ani}^{pl} = \psi_{kin}^{pl} - \psi_{dis}^{pl} \quad (20)$$

where ψ_{iso}^{pl} , ψ_{ani}^{pl} , ψ_{kin}^{pl} , and ψ_{dis}^{pl} represent the constituents of plastic free energy density from isotropic, anisotropic, kinematic, and distortional deformations, respectively (ψ_{dis}^{pl} is related to the distortional strain hardening with directional distortion of the yield surface but exploring this relation is not in the scope of this work). We can calculate ψ_{iso}^{pl} and ψ_{kin}^{pl} via [55]:

$$\psi_{iso}^{pl} = \frac{c_1}{2\rho} \bar{k}^2; \quad \psi_{kin}^{pl} = \frac{c_2}{2\rho} \bar{\alpha}_{ij} \bar{\alpha}_{ij} \quad (21)$$

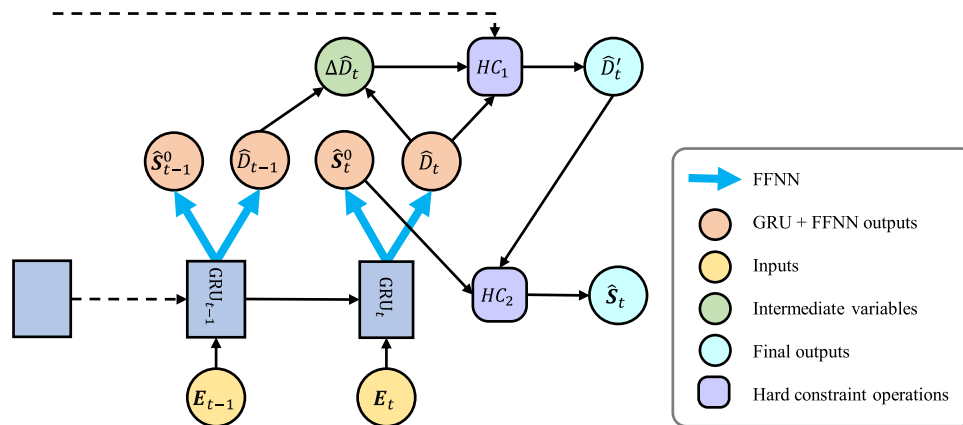


Fig. 4 Computational graph of physics-constrained RNN architecture: In this architecture, we illustrate the GRU cells at different time instances with their corresponding inputs and outputs. Intermediate variables are obtained from GRU+FFNN outputs and then used to impose hard constraints. HC_1 corrects \hat{D}_t to be non-decreasing, and

HC_2 computes the damaged stress \hat{S}_t using the undamaged stress \hat{S}_t^0 and the corrected effective damage parameter \hat{D}'_t . In addition, we use solid arrow lines to represent the data flow directly associated with the last two time steps shown in this figure and dash lines to represent the data flow from the previous time steps

where \bar{k} and $\bar{\alpha}_{ij}$ are, respectively, the thermodynamic conjugates to the size of the yield surface and the deviatoric back stress tensor that represents the center of the yield surface, and ρ is the material density. c_1 and c_2 are two non-negative material constants depending on the type of material models. We can therefore express the total plastic free energy as:

$$W^{pf} = \int_{\Omega} \left(\frac{c_1}{2\rho} \bar{k}^2 + \frac{c_2}{2\rho} \bar{\alpha}_{ij} \bar{\alpha}_{ij} \right) d\Omega_{0m} \geq 0 \quad (22)$$

By plugging Equations (16, 18, 22) into Equation (14), we show that for an RVE (associated with a macro element) that is subject to a generic deformation with hardening and softening, its effective macroscale stress and strain satisfy the following constraint:

$$\int_{\Omega} S_M dE_M \geq 0 \quad (23)$$

3.3.2 Hard Constraint: Damage Parameter

Our RNN is designed to predict the current stress tensor (S_t) and damage variable (D_t) given the history of the strain tensor (E_0, \dots, E_t). That is, a sequence of strain tensors (up to and including the current load step) is the input of the RNN while the current stress tensor and damage parameter are the RNN's outputs.

As schematically shown in Figure 4 we make several important changes to the vanilla data-driven RNN of Sect. 3.2 which directly learns the relation between (E_0, \dots, E_t) and (S_t, D_t). Specifically, we append the outputs of the RNN cells with two FFNNs and choose (S_t^0, D_t) as the outputs, that is, we use the effective reference stresses (which exclude

the effect of damage) as opposed to the damaged stresses. In addition, by considering the fact that our studied material is not self-healing during deformation, we hard-code the physical requirement on the damage parameter into the network such that its predictions are always non-negative and non-decreasing for any arbitrary deformation path that is fed to the network.

The rationale behind learning S_t^0 instead of S_t is that it forces the state variables of our RNN to learn hardening (characterized via the effective reference stresses) and softening (characterized by the damage parameter) in a decoupled manner. The damaged stresses are then obtained by combining the two output pairs in the FFNNs as (the hats indicate that the variables are predicted by the network):

$$\hat{S}_t = (1 - \hat{D}'_t) \hat{S}_t^0 \quad (24)$$

which is used by continuum damage mechanics, see Equation (A-3). The advantage of this decoupling or hard-coding the RNN to learn (S_t^0, D_t) is that it reduces the network size as the network no longer needs to internally learn the relation in Equation (24). Given that we cannot build a particularly large training dataset, small networks are preferred in our application.

Since our material does not heal itself during the irreversible damage process, \hat{D}_t should be non-decreasing along any deformation path. We mathematically represent this requirement as:

$$\dot{D}_t = \frac{\partial D_t}{\partial t} \geq 0 \quad (25)$$

where \dot{D}_t is the damage rate, and D_t is the effective damage parameter at time step t that can be computed from an RVE's effective damaged and reference stresses, see Equation (A-4). Equation (25) is not necessarily satisfied by the vanilla data-driven network in Fig. 3a, and thus we develop an efficient numerical scheme to explicitly enforce it. To this end, we first use \hat{D}_t and \hat{D}_{t-1} to compute the damage increment $\Delta \hat{D}_t$ and then update the damage parameter via the following scheme that ensure the damage increments are always non-negative:

$$\hat{D}'_t = \hat{D}_t + \sum_{\tau=1}^t \left(\Delta \hat{D}_\tau \times \left(0.5 \times \text{sign}(\Delta \hat{D}_\tau) - 0.5 \right) \right) \quad \forall t \in \{1, 2, \dots, n_{load} - 1\} \quad (26)$$

where we note that $\Delta \hat{D}_\tau = \hat{D}_\tau - \hat{D}_{\tau-1}$ and $\hat{D}'_0 = \hat{D}_0 = 0$. \hat{D}'_t indicates the corrected effective damage parameter which is then used in Equation (24) to compute the damaged stresses. The predicted (\hat{S}_t, \hat{D}'_t) are then compared to the ground truth as described in Equation (28) to minimize the loss function as described in Sect. 3.3.3.

3.3.3 Formulation of Loss Function

We design a composite loss function to be minimized via mini-batch stochastic gradient descent. The first component of this loss is the reconstruction error which at the arbitrary time instance $t \in \{1, 2, \dots, n_{load}\}$ is calculated as:

$$l_t^0 = \frac{1}{d_{out}} \frac{1}{n_b} \sum_{b=1}^{n_b} \| \mathbf{y}_t^b - \hat{\mathbf{y}}_t^b \|_2 \quad (27)$$

where d_{out} is the dimension of outputs, $\| \cdot \|_2$ indicates the l^2 norm of vectors, n_b is the size of the mini-batch, and \mathbf{y}_t and $\hat{\mathbf{y}}_t$ represent the ground truth and predicted values, respectively. In our case, $\mathbf{y}_t = (\mathbf{S}_t, D_t)$ and hence $d_{out} = 7$.

The second part of our loss function follows Equation (23) of Sect. 3.3.1 which indicates that the total internal work at an arbitrary macro IP can be computed from its associated RVE's homogenized stress and strain tensors, and that its value should be non-negative at any time instance:

$$l_t^1 = \frac{1}{n_b} \sum_{b=1}^{n_b} ReLU \left(- \sum_t \left(\hat{\mathbf{S}}_t^b : \Delta \mathbf{E}_t^b \right) \right) \quad (28)$$

where we approximate the total internal work at time step t by summing incremental internal works (which are computed by the predicted current stress $\hat{\mathbf{S}}_t^b$) and the incremental strain $\Delta \mathbf{E}_t^b$ in a training batch, i.e., $\Delta \mathbf{E}_t^b = \mathbf{E}_t^b - \mathbf{E}_{t-1}^b$. We use the rectified linear unit (ReLU) defined as $ReLU(x) = \max(x, 0) \geq 0$ in Eq. 28 to ensure the total loss is not penalized if the internal work is positive. We define the RNN's

composite loss function as:

$$\mathcal{L} = \sum_{t=1}^{n_{load}} l_t; \quad l_t = l_t^0 + \lambda l_t^1 \quad (29)$$

where λ is a scalar that controls the contributions of l_t^1 to the overall loss. To estimate an appropriate value for λ , we train a few RNNs and choose the one whose corresponding RNN achieves the smallest error on validation data. In this work, we use $\lambda = 10^{-6}$. A more systematic way to determine λ is through an adaptive weighting study [56].

As described above we enforce the two physics constraints within our RNN architecture by using two different approaches. While we implement the energy constraint in Eq. 28 as a soft constraint by adding an associated penalty term in the loss function, we enforce the damage constraint in Equation (25) as a hard constraint by imposing architectural modifications and post-processing RNN cells' outputs by using the intermediate variables in the network. While the hard constraint always guarantees that the required conditions are met, it may lead to a stiffer optimization problem [57] in training. Additionally, hard constraints require architectural modifications which may be difficult to implement for complex physical conditions. In our studies, the architecture in Fig. 4 resulted in more accurate models compared to cases where the constraint in Equation (25) was enforced by penalizing the loss function.

3.3.4 Teacher Forcing

The networks presented so far have recurrent connections between the output at one time step and the hidden units at the next time step. Compared to models with hidden-to-hidden connections, these networks are less powerful since their outputs should capture all of the information about the past that is needed in predicting the current state. One technique for addressing this issue is teacher forcing [58] which refers to networks whose outputs are fed back into the model via recurrent connections. Figure 5 schematically demonstrates such a model with one look-back step where only the previous output is fed back into the model (more look-backs are possible).

As illustrated in Fig. 5, training and testing are done slightly differently in the presence of teacher forcing. At the offline training stage, we provide the ground truth (RVE's effective stresses and damage variable) at the previous time step as inputs at the next time step. At the online testing stage, we must use the predictions at the previous time step since the ground truth is unavailable.

Teacher forcing reduces the training time and provides smaller closed-loop reconstruction errors. However, networks with teacher forcing typically struggle in open-loop

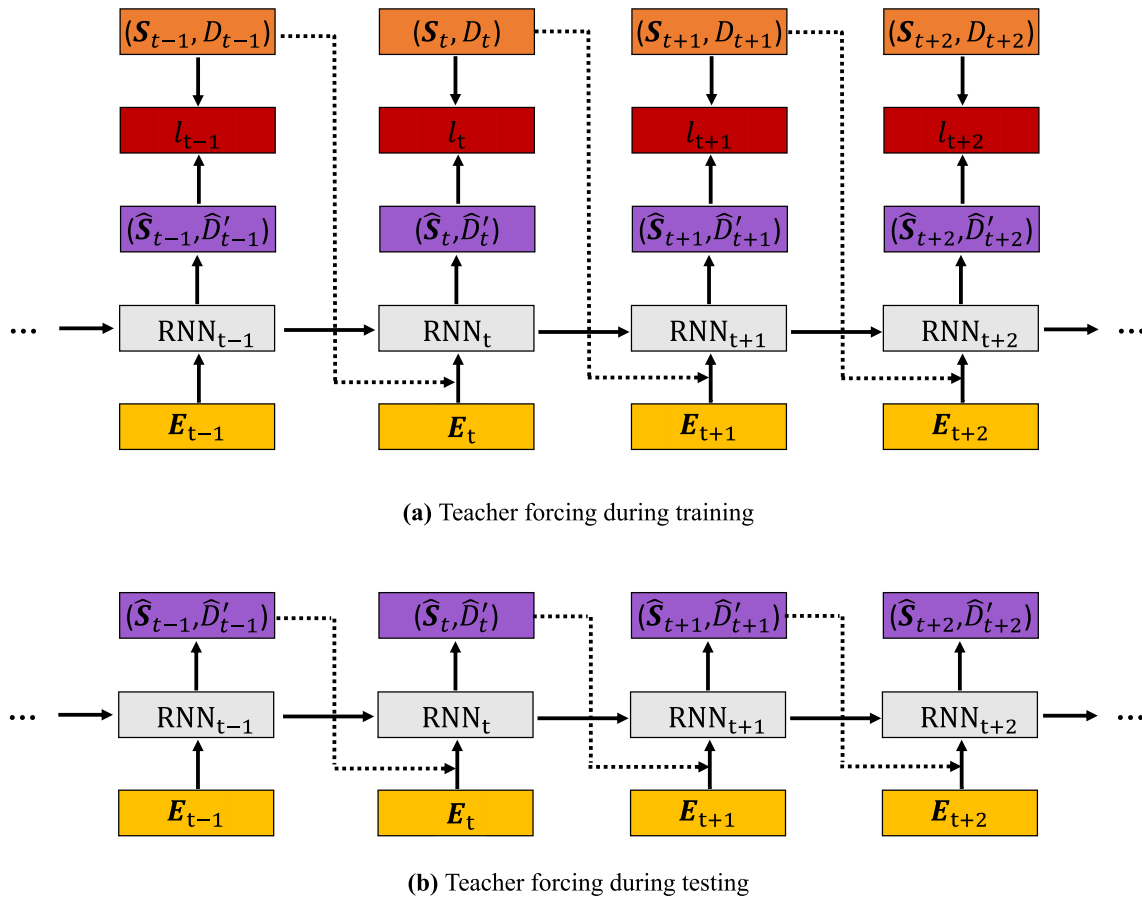


Fig. 5 **a** In training, the current input is augmented with the ground truth from the last time step. **b** During testing, the inputs are augmented with the predictions from the previous time step

applications where the ground truth is unavailable. While some training techniques (such as intentionally corrupting the outputs before they are fed back into the model) alleviate this issue to some extent, large errors can still appear. As we demonstrate in Sect. 4, using an RNN in a multiscale simulation is indeed an open-loop application where teacher forcing provides poor performance.

3.4 Surrogate Integration with Multiscale Solvers

Our trained RNN surrogates the computationally expensive micro-solver in a multiscale simulation. The online deployment of the RNN surrogate within an iterative solver poses some difficulties. During training the RNN has access to the deformation and effective response histories at all n_{load} load steps. Comparatively, in the online computations, only the previous strains and responses are available and even the converged macro strains at the current load step are unknown (in conventional numerical solvers these macro strains are found iteratively by solving the equilibrium equations, see Sect. 2). We address these difficulties by explicitly modifying

the input sequences and implicitly resetting RNN's hidden variables amid iterations.

We embed our trained RNN in a multiscale model per the pseudo-code in Algorithm 1 which primarily aims to integrate the RNN with the Newton–Raphson method. In the presence of nonlinear deformations, the Newton–Raphson method is typically used to iteratively solve for the material's path-dependent responses. This method essentially consists of a double-loop structure where the outer loop incrementally steps along the applied load while the inner loop aims to match internal and external forces by iteratively updating the material response under a fixed loading condition.

In a typical step of a multiscale simulation, we compute the macro strain tensor at an arbitrary IP from equilibrium equations within the inner loop. By appending the strains at the current iteration to the sequence of previous convergent strains, the length of the strain sequence equals the current load step number that is smaller than n_{load} (n_{load} is the required length of the RNN's input sequence). To address this mismatch, we repeat or pad the value of the current strain multiple times at the end of the strain sequence. This padding only makes the strain sequence compatible with the RNN's

```

i = 1, 2, ..., nload; /* Newton step number
*/
j = 1, 2, ..., niter; /* Newton iteration
number */
k = 1, 2, ..., nmip; /* Number of macroscale
IPs */
 $\epsilon = 10^{-6}$ ; /* Convergence criterion */
while i ≤ nload do
  while j ≤ niter do
    (1) Read macro strain  $\mathbf{E}_i^j$  from macro
    equilibrium equation
    (2) Append  $\mathbf{E}_i^j$  to the convergent strain sequence
     $\{\mathbf{E}_1^c, \mathbf{E}_2^c, \dots, \mathbf{E}_{i-1}^c, \mathbf{E}_i^j\}$ 
    (3) Add (nload - i) replicate padding of  $\mathbf{E}_i^j$  to the
    end of the sequence in (2)
    while k ≤ nmip do
      (4) Perform RNN inference on the updated
      strain sequence for each macro IP
    end
    (5) Retrieve RNN's outputs for the effective
    responses at the step i
    (6) Solve macro equilibrium equation
    if  $\|\mathbf{f}_{int} - \mathbf{f}_{ext}\| \leq \epsilon$  then
      Update convergent strain  $\mathbf{E}_i^c = \mathbf{E}_i^j$ 
      Continue to the next load step: i ← i + 1
      Break; /* Iteration convergence
      */
    else
      Continue to the next iteration: j ← j + 1
    end
  end
end

```

Algorithm 1: Integration of RNN in multiscale analysis

input, but also implicitly freezes the RNN's hidden variables at the current step within the iterations (inner loop). This freezing happens because the values of the hidden variables at the current time instance are decided by the state of network parameters and the inputs from previous time instances, see Equation (9). This freezing is akin to the classic radial return algorithm in plasticity models where material state variables are only updated upon convergence. We also emphasize that the number of load steps in the multiscale simulation should be smaller than or equal to the sequence length of RNN inputs (n_{load}), since a larger step number results in data truncation during input data preparation and erroneous RNN inference (in our case, n_{load} is chosen large enough to prevent truncation).

4 Numerical Experiments

In this section, we first illustrate the efficacy of our physics-informed RNN in predicting microstructural effective behaviors in Sect. 4.1 where we assume micro porosity as the only material defect. We then perform the computation of multiscale elasto-plastic hardening and softening simulations in Sect. 4.2 by integrating our RNN (as a surrogate of microstructural analysis model) with a macro FEM solver. In Sect. 4.3, we deploy our multiscale surrogate model to perform a mesh convergence study on a component with different spatial discretization levels to simulate macro damage patterns. In the experiments, we record computational costs and perform accuracy analysis.

We build our RNN via TensorFlow in Python. We generate the database of microstructural effective responses on a state-of-the-art high-performance cluster (HPC) which has 60 CPU cores (AMD EPYC processors) and 360 GB RAM. We train the RNN model on the HPC via two GPU units (NVIDIA Tesla v100) with 32 GB RAM. For multiscale simulations, we develop a dedicated program to integrate our RNN model within a multiscale analysis engine which is implemented in Matlab. We note that all data-driven multiscale computations in Sects. 4.2 and 4.3 are conducted on a 64-bit Windows desktop with four CPU cores (Intel i7-3770) and 16 GB RAM.

4.1 Surrogate of Microscale Damage Modeling

We aim to use the proposed surrogate to accelerate damage analysis for metallic alloys with process-induced porosity. Specifically, we analyze cast aluminum alloy A356 whose heterogeneity is represented by microstructures with two material phases: a void phase (representing porosity) and a material phase. The latter phase is presumed to have isotropic homogeneous properties (without considering local material features such as impurities or grain structures and interfaces). Therefore, we can simulate damage propagation within the material phase by using classic continuum damage mechanics as argued in Section A.

We assume that the A356 has an elastic modulus of 5.7e4 MPa and a Poisson's ratio of 0.33, and its isotropic elasto-plastic hardening behavior follows an associative plastic flow rule with the Von Mises yield surface defined by:

$$S \leq S_y(\bar{E}^{pl}) \quad (30)$$

where S and S_y are, respectively, the Mises equivalent stress and the yield stress which depends on the equivalent plastic strain \bar{E}^{pl} . To model strain hardening, we assume the relation between S_y and \bar{E}^{pl} is piecewise-linear as shown by the hardening curve in Fig. 6. For modeling softening, we employ the damage continuum model discussed in Section A

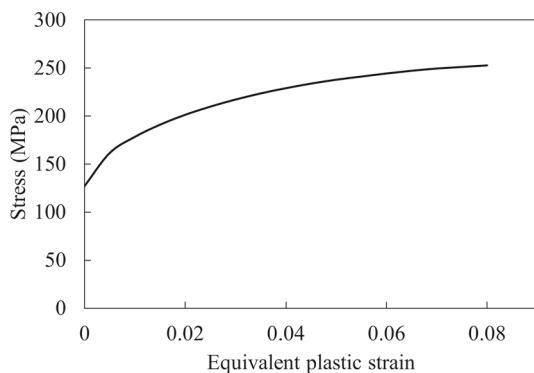


Fig. 6 Elasto-plastic hardening behavior: We use a piecewise linear hardening model to define the elasto-plastic behavior of aluminum alloy A356

with the fracture strain E_f of 0.067 and the fracture energy G_f of $1.92e4$ N/m.

4.1.1 Database Generation of Microstructural Effective Responses

We solve micro BVPs on a microstructure whose geometry and mesh are illustrated in Fig. 7a. The microstructure is made out of A356 and contains a spherical pore at its center. Even though classic FEM with sufficiently fine mesh, e.g., see Fig. 7b, can provide high-fidelity solutions to BVPs, it is generally expensive, especially for the response database generation. To improve computational efficiency, we apply mechanistic DCA-based ROM [25, 26, 59] to solve the BVPs. Compared to FEM, our ROM strikes a good balance between efficiency and accuracy by agglomerating elements into clusters, e.g., see Fig. 7c where material IPs in the same cluster are assumed to share identical elasto-plastic behaviors.

We note that a mesh independence study is often required in material softening simulations to choose a proper spatial discretization for solution convergence [59, 60]. We conduct the microscale mesh convergence investigation in Appendix C where we systematically compare the softening behaviors between FEM and ROM for the RVE in Fig. 7a and show that our ROM with 1, 200 clusters predicts consistent post-failure behaviors as the FEM while being 10 times faster. Hence, we choose the ROM with 1, 200 clusters to build our database and consider this ROM as the ground truth when validating the data-driven surrogates in the following experiments.

To generate the database, we set the sampling constraint for the strain magnitude as $\zeta_1 = 10\%$ and the constraint of the volumetric strain as $\zeta_2 = 4\%$. For the GP interpolation, we set the number of control points with random strain values as $n_c = 5$. Our database contains a total to $n_p = 30,000$ deformation paths and RVE effective responses where each path includes six strain components, six effective stress com-

ponents, and one effective damage variables at $n_{load} = 101$ sequential loading steps. Generating this database costs about ten-day computational time on the HPC by exploiting parallel computing with 60 CPU cores.

4.1.2 Impacts of Physics Constraints

To demonstrate the impacts of the two physics constraints in Sect. 3.3, we compare the prediction accuracy of a pure data-driven vanilla model against our RNN model. For this comparison, we randomly choose 200 deformation-responses sequences as a test set. We further randomly select 6,000, 12,000, 18,000, 24,000 and 29,800 sequences from the database to form five different training-validation datasets. For all training-validation datasets, we split them into 80% for training set and 20% for validation set. For example, the dataset of the size of 6,000 has 4,800 sequences for training and 1,200 for validation. We point out the training and validation sets serve different purposes, as the training set is used to iteratively update learning parameters during BPTT, while the validation set is used to detect overfitting or underfitting.

During training, we normalize all data sequences and use 1,200 epochs with a batch size of $n_b = 64$. We choose *Adam* as the optimizer with an adaptive learning rate that starts at 10^{-3} and reduces by 25% when the validation error is not reduced over 30 training epochs. We terminate the training process when the training reaches the maximum number of epochs or the loss function is not improved by 10^{-7} over 50 epochs. We use mean squared error (MSE) to measure accuracy:

$$\text{MSE} = \frac{1}{n_t n_{load} d_{out}} \sum_{m=1}^{n_t} \sum_{t=1}^{n_{load}} \sum_{i=1}^{d_{out}} (y_{i,t}^m - \hat{y}_{i,t}^m)^2 \quad (31a)$$

$$\text{MSE}_S = \frac{1}{n_t n_{load} d_S} \sum_{m=1}^{n_t} \sum_{t=1}^{n_{load}} \sum_{i=1}^{d_S} (S_{i,t}^m - \hat{S}_{i,t}^m)^2 \quad (31b)$$

$$\text{MSE}_D = \frac{1}{n_t n_{load}} \sum_{m=1}^{n_t} \sum_{t=1}^{n_{load}} (D_t^m - \hat{D}_t^m)^2 \quad (31c)$$

where MSE accounts for the total prediction error including both stress and damage predictions, while MSE_S and MSE_D are the prediction error for stress and damage, respectively. n_t and d_{out} are the number of data sequences in the test set and the dimension of outputs. $y_{i,t}^m$ and $\hat{y}_{i,t}^m$ are the ground truth and prediction for the i^{th} output component at the t^{th} load step for the m^{th} test sample. In addition, d_S , $S_{i,t}^m$, $\hat{S}_{i,t}^m$, D_t^m and \hat{D}_t^m are the number of 3D stress components, true stress, predicted stress, true damage and predicted damage variables, respectively, i.e., $y_{i,t}^m = (S_{i,t}^m, D_t^m)$. We note that the values of n_t , d_{out} and d_S are equal to 200, 7 and 6, respectively.

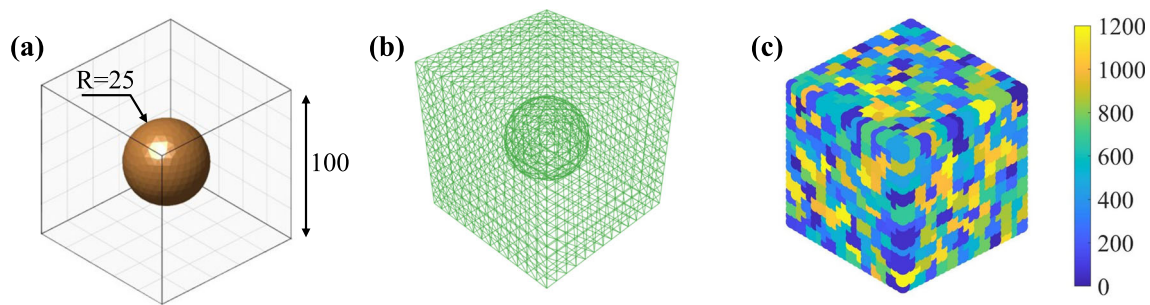


Fig. 7 The geometry, dimension and mesh of our RVE: **a** The dimension (unit: μm) of our RVE that contains a spherical pore at center with a pore volume fraction of 6.25%; **b** The RVE is discretized by 15, 000 finite elements; and **c** Our ROM with 1, 200 clusters

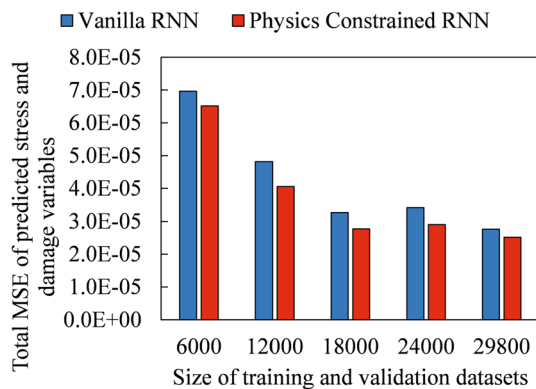


Fig. 8 Convergence study: We compare the total MSE between the pure data-driven model and the proposed physics constrained RNN model to demonstrate the effectiveness of using constraints and domain knowledge in designing the RNN. The errors are computed on the same test dataset which has 200 deformation-response sequences

Furthermore, we set the penalty parameter as $\lambda = 10^{-6}$ in the customized loss function in Eq. 29.

After the 10 models are trained on the five different training-validation datasets, we compare their prediction errors on the same test set that is unseen amid the entire training process, see Figs. 8 and 9. For the overall MSE, we find that as the sizes of training-validation datasets increase from 6, 000 to 29, 800, the MSEs of both models decrease dramatically from about 7×10^{-5} to 3×10^{-5} . It is clear that the overall MSE of our proposed model is always lower than that of the vanilla model.

To provide more insights into the performance of two models across the different data set sizes, we report the MSEs that each model achieves in terms of stresses and damage variable, see Fig. 9 and Equation (31). It is evident that the proposed model provides a higher accuracy than the vanilla model. Particularly, on the smaller datasets (6, 000 or 12, 000 sequences) with a limited amount of training data, the prediction error of our physics constrained RNN is about 60% lower than the vanilla model. As the size of the training data increases, we observe that the gap decreases. A similar trend

is observed for the damage variable except that the gap in the performance is smaller.

To visualize the predicted effective responses by our physics-constrained surrogate model, we randomly select four strain paths from the test set and compare our predictions against the ground truth, see Fig. 10. We observe that our RNN provides accurate predictions for all effective stresses and the damage variable even though the deformation paths are quite complex. In particular, we observe that as the damage variable increases to 1.0 amid material deformation, the magnitudes of the effective stresses are correspondingly reduced (which indicates that the RVE's load-carrying capacity significantly decreases).

Our surrogate model approximates the elastoplastic-damage behaviors of porous metallic alloys which typically fracture within small strain ranges. As shown in Fig. 10, the typical strain ranges of our studied material in realistic applications are relatively small, i.e., the metallic aluminum matrix fractures before high strains occurs. Hence, we can assume the results based on the small and finite strain formulations are similar. In our dataset generation step of Sect. 3.1, we went above this range to ensure that we sufficiently sample deformation paths with different strain values. For the above reasons, we presume that an infinitesimal strain theory can be used for simulating the deformation of A356. A more accurate approach can be obtained by leveraging a finite strain theory in the data generation step.

4.1.3 Impacts of Teacher Forcing

As discussed in Sect. 3.3, teacher forcing, which augments ground truth or predictions from previous steps to the input at the current step during training, can provide an effective way for improving the accuracy of RNNs. To quantify the impact of teacher forcing, we compare the total MSE of the predicted stress and damage variables over the test dataset between our physics-constrained RNN model with and without the teacher forcing technique in Table 1.

We implement two teacher forcing models here: the first model with the number of look back step (NLB) of one, and

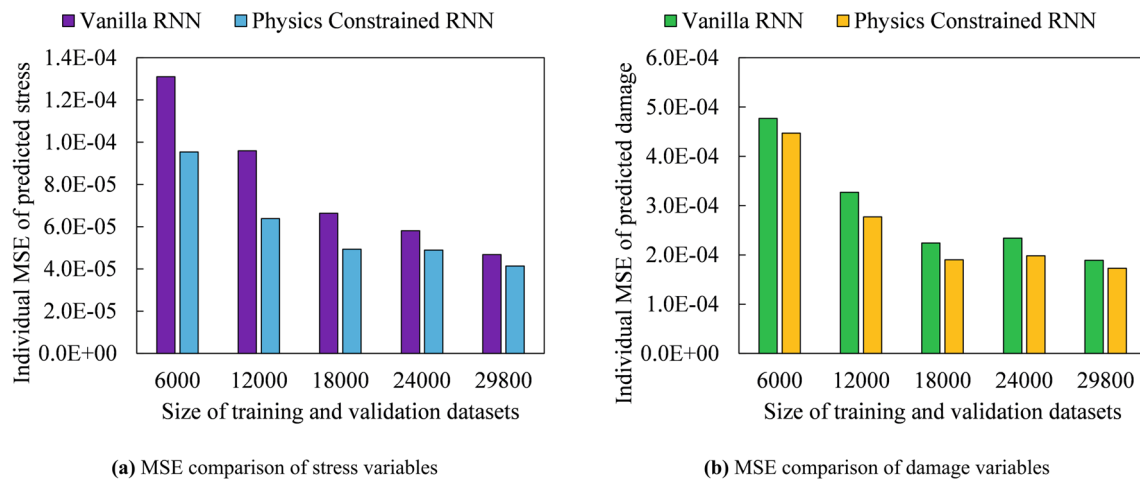


Fig. 9 MSE on predicted stress and damage: We compare the vanilla RNN and our physics constrained RNN models based on MSE in stresses and damage variables

the second model with the NLB of five. From Table 1, we observe that compared to the model without teacher forcing, the two teacher forcing models reduce the total MSE by 21.4% (NLB=1) and 24.2% (NLB=5), respectively. Comparing the individual MSEs, we find that teacher forcing reduces the prediction error of effective damage while not for the stress. Therefore, in single-scale RVE simulations, the teacher forcing improves our RNN's overall prediction accuracy. However, as we show next, teacher forcing significantly reduces the performance in multiscale simulations.

4.2 Surrogate of Multiscale Damage Modeling

After we demonstrate that our RNN can accurately predict microstructural effective responses under various deformation paths in Sect. 4.1. We can now use the RNN as a faithful surrogate to replace the computationally expensive microstructural analysis in multiscale simulations.

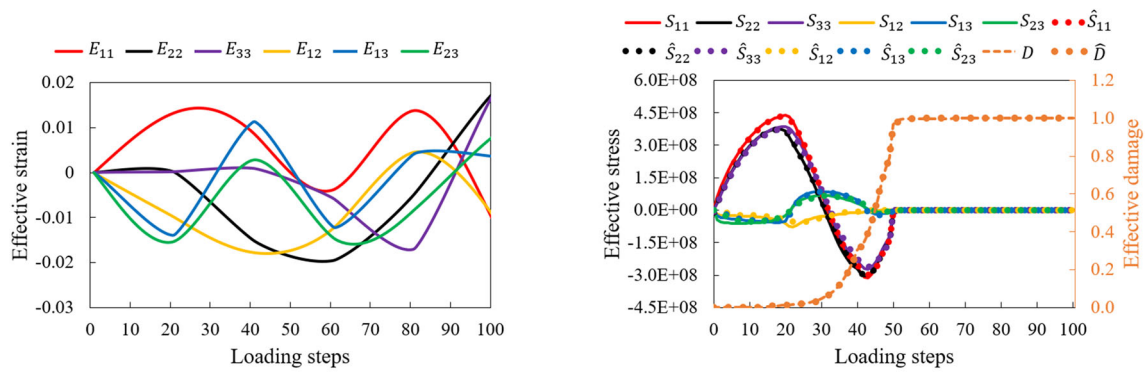
We design a 3D L-shape bracket for our multiscale simulation as in Fig. 11a. The bracket is subject to a Dirichlet boundary condition on the left side along the x-axis while its right surface is fully fixed in all directions. Based on single-scale damage simulations on the same L-shaped bracket [59] it is observed that strain concentrations appear around its sharp corner and result in damage propagation from the corner to the outer surface. Therefore, we create a sufficiently large multiscale region around this zone to associate macro elements with porous microstructures. Specifically, we associate each IP of the multiscale domain with a porous RVE as illustrated in Fig. 7. To save computational costs, we assume there is a mono-scale domain outside the multiscale domain where the IPs are not associated with any RVEs. We mesh the bracket by 5, 300 tetrahedral elements of reduced integration. The multiscale domain contains 360 elements each

of which is associated with an RVE that is decomposed by 1, 200 clusters.

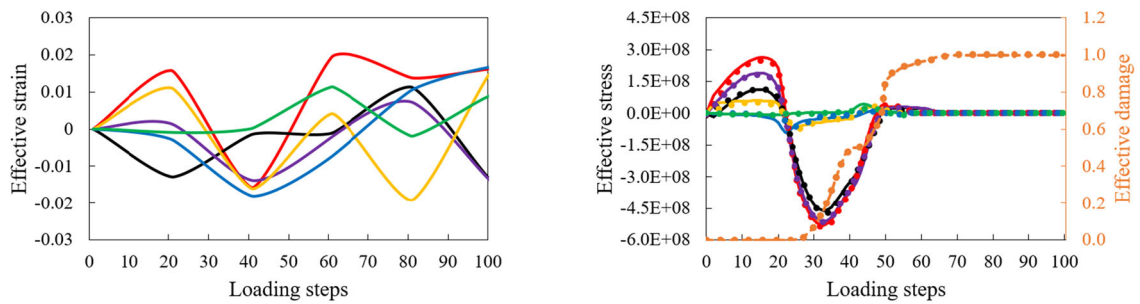
We first demonstrate the accuracy of the multiscale model for elasto-plastic hardening behaviors under complex cyclic loading histories. To this end, we subject the bracket to a loading-unloading-reloading condition by setting the Dirichlet boundary condition as $d = 0 \rightarrow 2 \rightarrow 0 \rightarrow -2$ mm. We compare the resulting reaction force and displacement curves between our proposed FE-RNN approach and the benchmark FE-ROM method in Fig. 11b.

As shown in Fig. 11b we observe that teacher forcing, either with one or five look back steps, provides erroneous results which are due to the fact teacher forcing leverages the historical predictions (i.e., stresses and damage variable predicted for previous load steps) in estimating the current stresses and damage variables. However, these historical predictions are highly noisy since they suffer from errors that are (1) accumulated: since erroneous predictions are fed back into the model and propagated forward (see Fig. 5), and (2) compound: since the predictions at any macro IP affect the predictions at other IPs. Following these results, we adopt our FE-RNN model without teacher forcing for all multiscale simulations in the following experiments (note that this multiscale simulation is different from the single-scale study in Sect. 4.1.3 where we use different RNNs to surrogate the effective responses of an RVE).

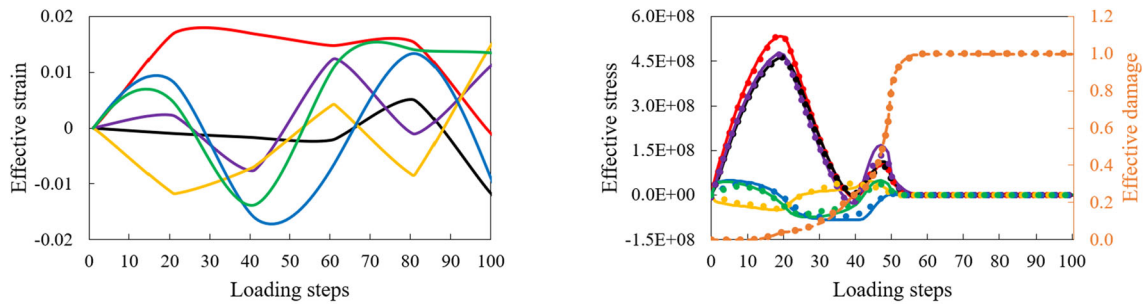
We compare the Von-Mises stress distributions between the benchmark and our FE-RNN model by setting the boundary condition as $d = -2$ mm, see Fig. 12. We observe a good agreement between the two models and only observe minor local discrepancies at the sharp corner. These errors are primarily due to the fact that RNN's prediction accuracy decreases at extreme cases which are insufficiently sampled at the training stage. We also note that deep NNs (including



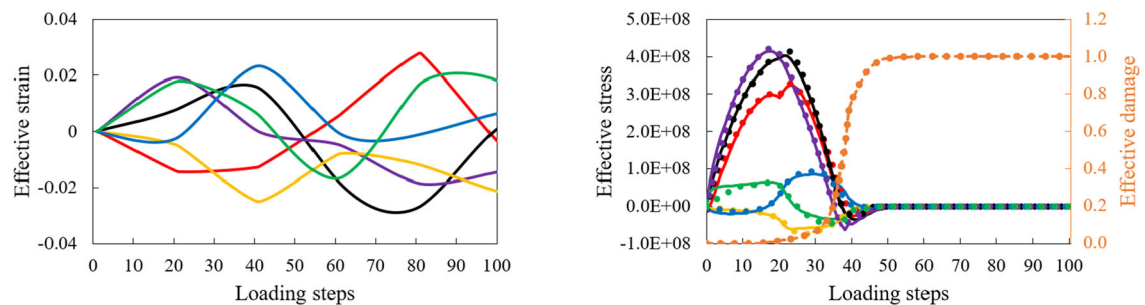
(a) Random strains and effective responses of test sample 1



(b) Random strains and effective responses of test sample 2



(c) Random strains and effective responses of test sample 3



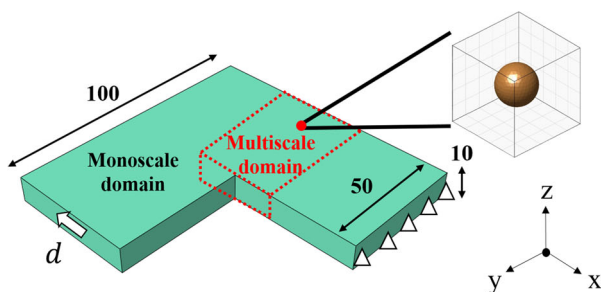
(d) Random strains and effective responses of test sample 4

Fig. 10 RVE deformation and responses from RNN versus ground truth: We demonstrate four test samples with different strain paths (the first column where each path contains six strain components with 100 steps) and their associated effective stresses and damage variable (the second column)

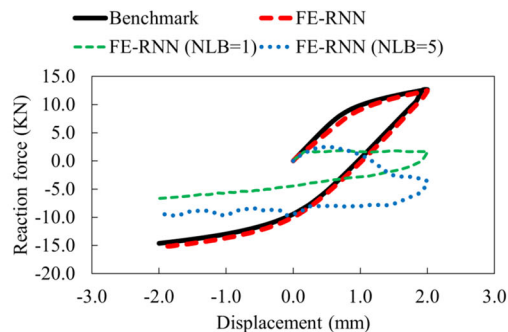
Table 1 Impacts of the teacher forcing on single-scale predictions: Comparison of the total MSE and individual MSEs (10^{-5}) of the predicted effective stress and damage variables over the test set between our

physics constrained RNN model (trained by 29,800 sequences) without teacher forcing and the same model with teacher forcing using one or five look back steps

	No teacher forcing	Teacher forcing (NLB=1)	Teacher forcing (NLB=5)
Total MSE	2.52	1.98	1.91
MSE _S	4.14	4.62	4.54
MSE _D	17.30	9.67	9.24



(a) Geometry, dimensions (unit:mm) and boundary conditions



(b) Reaction force and displacement

Fig. 11 Multiscale model of the L-shape bracket: **a** Every macroscale IP in the multiscale domain is associated with a microscale porous RVE; and **b** Comparison of the load–displacement curves of the elasto-plastic

hardening behaviors between benchmark and the proposed FE-RNN models with and without teacher forcing

RNNs) are notorious for poor extrapolation ability and provide poor predictions for rare events that are insufficiently seen during training (e.g., the sharp drop in the stress–strain curve upon fracture).

In our second multiscale experiment, we simulate the elasto-plastic hardening and softening of the same L-shape bracket where its Dirichlet boundary condition is set as $d = 10$ mm. To prevent the occurrence of the non-physical single-layer fracture bands as discussed in Section A, we apply a non-local damage function (see Eq. A-5) with a strain localization bandwidth of $l_d = 15$ mm, see Fig. 13a for a comparison between l_d and the mesh size of the bracket. The force-displacement curves are compared in 13b where the general trends of the two methods match well especially in the hardening section. Minor discrepancy manifests in the softening regime where the data-driven model tends to break earlier which underestimates the component’s load-carrying capacity by about 2.5%. The underlying reason is that softening behaviors dramatically increase the complexity of the material’s governing equations and our training data has much less information on softening than on hardening.

We now compare the distributions of damage variables and Von-Mises stresses when the boundary condition is set

as $d = 10$ mm, see Figs. 14 and 15, respectively. We see that both field variables have good agreements between the two approaches. In Fig. 14, we observe that the fracture bands initiate from the sharp corner and stretch towards the right surface. We can also clearly see the effects of imposing non-local damage functions in avoiding non-physical single-layer fracture bands. As for the stress distributions in Fig. 15, both approaches show that the local stress values are significantly reduced within fracture bands that indicates a loss of load-carrying capacity in the fractured elements. We observe a minor discrepancy of local stresses at the front tip of the fracture bands between the two methods: while the benchmark indicates relatively low stresses at the highlighted region, our FE-RNN model suggests stress concentrations triggering more damage if the component is further deformed. These stress concentrations explain why our data-driven model predicts an earlier damage occurrence than the benchmark in Fig. 13b.

The discrepancy between the proposed model and benchmark can be further quantified by the histogram of errors as shown in Fig. 16. In terms of damage variables, it is quite clear from Fig. 16a that the two approaches yield identical solutions in the majority (more than 80%) of elements. Based

Fig. 12 Comparison of Von-Mises stress distributions in hardening simulation: **a** The ground truth distribution of Von-Mises stresses (unit: Pa) and **d** Von-Mises stresses by the proposed FE-RNN model

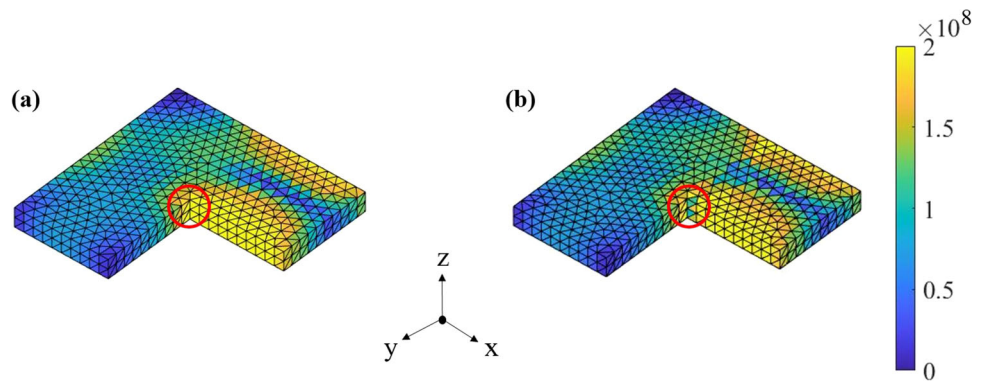


Fig. 13 Damage simulations of the L-shape bracket: **a** The macroscale discretization and strain localization bandwidth applied in the damage function; and **b** Comparison of the softening load–displacement curves between benchmark and the proposed FE-RNN model

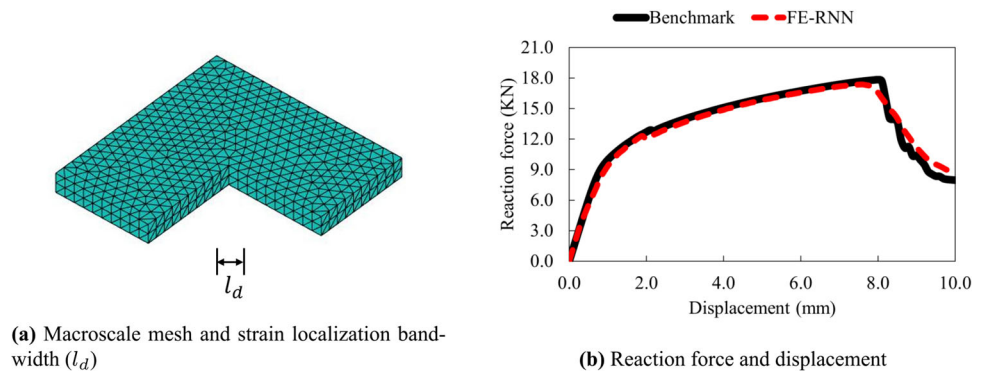


Fig. 14 Comparison of damage patterns: **a** The distribution of benchmark damage variables; and **d** Damage variables via our FE-RNN model where yellow indicates a full material fracture while blue represents an intact state

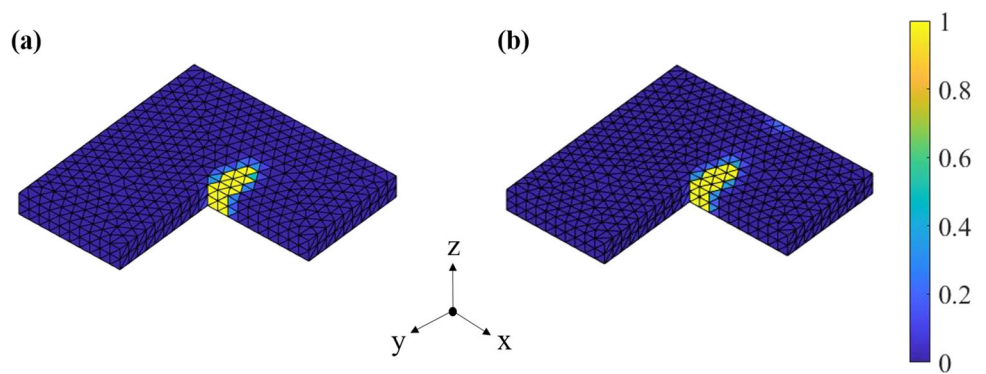
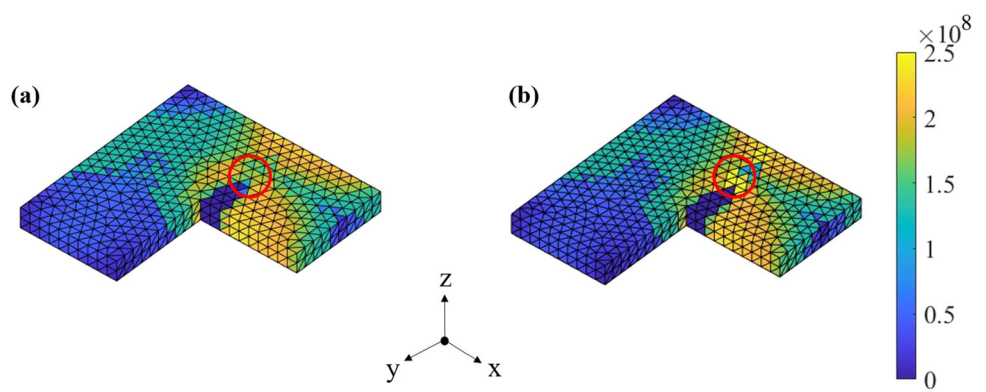


Fig. 15 Comparison of Von-Mises stress distributions in softening simulation: **a** The distribution of benchmark Von-Mises stresses (unit: Pa); and **d** Von-Mises stresses by our FE-RNN model



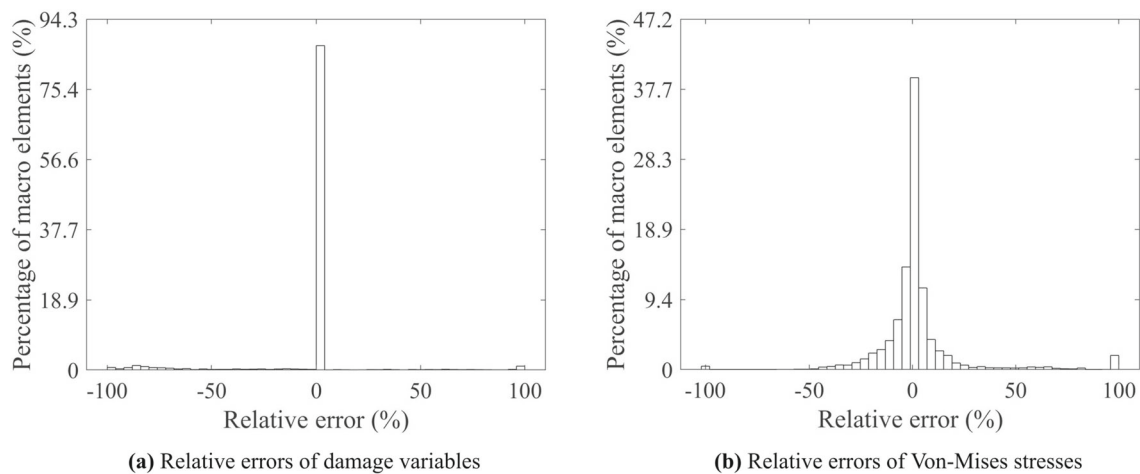


Fig. 16 Histogram of relative errors of field variables: **a** Relative errors of the values of damage variables between benchmark and our FE-RNN model in Fig. 14; and **b** Relative errors of the values of Von-Mises stresses in Fig. 15

Table 2 Breakdown of computational costs of the L-shape bracket model: Despite considerable costs in database generation and training of the RNN, its efficiency (measured by clock time) in the multiscale

simulations is about $125\times$ and $1240\times$ higher than ROM and FE^2 , respectively, where the time estimation of FE^2 comes from the time comparison in Fig. 21b

	FE-RNN (proposed)	FE-ROM (benchmark)	FE^2 (estimated)
Database generation	239.5×60 CPU-hour	-	-
RNN training	4.3×2 GPU-hour	-	-
Multiscale simulation	0.4×4 CPU-hour	49.8×60 CPU-hour	494.0×60 CPU-hour

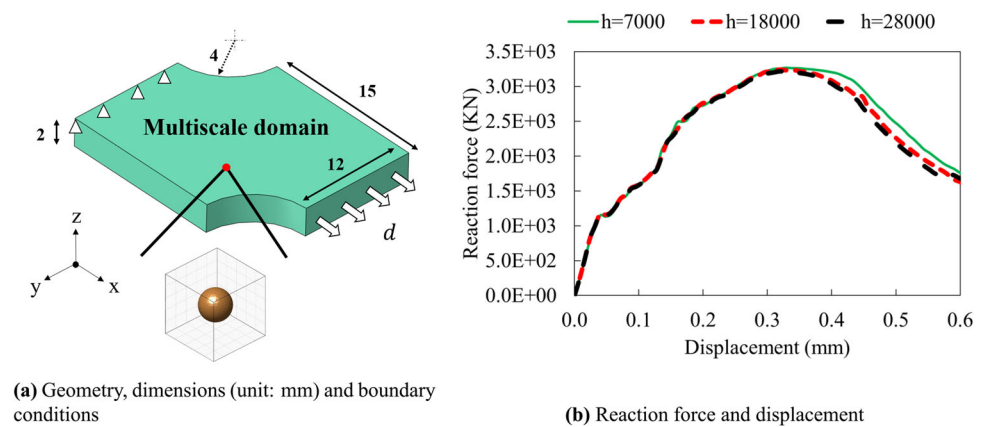
on the distributions of stress errors in Fig. 16b, we can see relatively large errors in about 2% of all elements. It should be noted, however, for most elements, their absolute errors are smaller than about 10% between the two methods.

To quantify computational efficiency, we break down the costs of different steps in this multiscale model as shown in Table 2. Compared to the mechanistic models (FE-ROM and FE^2), our data-driven model (FE-RNN) requires additional costs on database generation and model training. Even though expensive, we only need to perform the two steps once, and after training we can deploy the trained RNN model for any multiscale simulation without extra costs. In terms of the online clock time, our model shows superior efficiency to the mechanics models (FE-ROM and FE^2) with about $125\times$ and $1240\times$ accelerations, respectively. It is noted that we do not directly perform the FE^2 due to its prohibitively demanding costs, its computational time is estimated by comparing to the ROM on a smaller multiscale simulation whose time comparison is demonstrated in Fig. 21 of Appendix C. We also note that, while we perform the training process on two GPU processors, we carry out both the database generation of the FE-RNN and the multiscale simulations of FE-ROM and FE^2 by paralleling 60 CPU cores with 360 GB RAM on an HPC. Comparatively, our proposed RNN model only needs four CPU cores on a desktop computer for the multiscale

computation, providing feasible solutions to the engineers without accessibility to large computational resources.

We can observe from Table 2 that the generation of training database is rather computationally expensive. However, our surrogate model is still more advantageous than the classic FE^2 approach for several reasons: (1) Even though the offline database generation time (239.5×60 CPU hours) is much longer than the online simulation time of FE-RNN (0.4×4 CPU hours) and FE-ROM (49.8×60 CPU hours), it is still about 1.06 times faster than FE^2 (494.0×60 CPU hours). In fact, the reported time for FE^2 in Table 2 is an estimate since we cannot actually run such simulations due to excessive computer memory requirements which do not hamper our FE-RNN approach. (2) The training database is only generated once and the trained surrogate can be used in future without any additional offline costs. Hence, these costs are expected to be amortized with the repeated future uses of the model. For instance, for the mesh independence study in Sect. 4.3, the trained model is directly used and there is no need to generate a new data or train a new machine learning model. Additionally, the macrostructure in Sect. 4.3 is different from the macrostructure in Sect. 4.2 which again shows that we can use our trained data-driven constitutive model across multiple applications. This “transferability” feature does not exist in FE^2 : Even for the same macroscale

Fig. 17 Multiscale model of double-notched specimen: **a** Every macroscale integration point is associated with a microscale porous RVE; and **b** Convergence study of the softening load–displacement curves with different macro discretization levels



model, if the distribution of the porosity changes (e.g., some macroscale IPs have no pores while some other IPs are associated with a microstructure which has pores), FE² needs to be rerun which is very expensive. With our approach, we only need to rerun the online macroscale simulation whose cost is be much lower.

We point out that the main purpose of this L-shape bracket study is to verify the accuracy of our multiscale surrogate by comparing it to the FE-ROM benchmark (we do not intend to use a finer mesh to study mesh convergence of damage analysis in this example). As shown in Table 2, the computational time of the benchmark for this relatively coarse mesh is about 49.8×60 CPU hours (2.1 days) on a high-performance computer by paralleling 60 CPUs. This computational time would be significantly higher (exceeding our computing budget) if we choose a much finer mesh. For the mesh convergence study of damage modeling, we compare post-failure behaviors of multiscale damage models with much finer mesh discretization in Sect. 4.3.

4.3 Multiscale Damage Surrogate: Mesh Independence Study

One of the major challenges of using continuum mechanics to simulate softening is preventing the fracture bands from residing in single-element-wide layers. A popular solution to this challenge is to apply non-local functions to constrain damage patterns to different spatial discretization levels. To this end, we apply the proposed RNN model to a new 3D model in this section and assess its robustness in predicting damage behavior while changing the mesh size.

The geometry, dimensions, and boundary conditions of the double notched specimen are demonstrated in Fig. 17a. The left surface of the specimen is fixed and its right surface is extended by $d = 0.6$ mm along the x-axis. In this experiment, we model the entire specimen as a multiscale structure where each macro IP is associated with a porous RVE. For the mesh convergence study, we discretize the macro specimen

with three different mesh sizes: a coarse mesh with 7,000 elements, a medium mesh with 18,000 elements, and a fine mesh with 28,000 elements.

We demonstrate the reaction force–displacement curves in Fig. 17b which indicates that the three mesh levels achieve very close elasto-plastic hardening responses and are slightly different in the softening regime. Specifically, we note that as the mesh level increases from medium to fine, the post-failure force–displacement responses tend to converge.

We also probe the convergence by inspecting the stress distributions and damage patterns, see Fig. 18. In Fig. 18a and b, we can clearly see that at all mesh levels the damage initiates from the inner circular surfaces and propagates across the specimen as it deforms. The influence of imposing non-local function is evident: it not only successfully avoids non-physical single-element-wide damage layers, but also constrains the fracture bandwidth regardless of the mesh sizes. Additionally, as indicated in Fig. 18a stress concentrations consistently appear at both fracture front tips and around sharp corners.

We report the simulation time of this multiscale double notched specimen in Table 3. We emphasize that due to the superior efficiency, we can apply our trained FE-RNN to any multiscale models with no extra costs of data generation and model training. Additionally, our trained FE-RNN model is memory-light and can be run on a desktop with four CPU cores and 16GB RAM. Based on the time comparison in Table 2, simulation of the multiscale model with 28,000 elements requires a clock time of 1,304.8 h (54.4 days) and 12,942.8 h (539.3 days) by paralleling 60 CPU cores with 360 GB RAM by FE-ROM and FE² methods, respectively.

5 Conclusions

We propose a physics-constrained deep learning model that surrogates the homogenized path-dependent microstructural behaviors in 3D large-scale multiscale simulations. Our

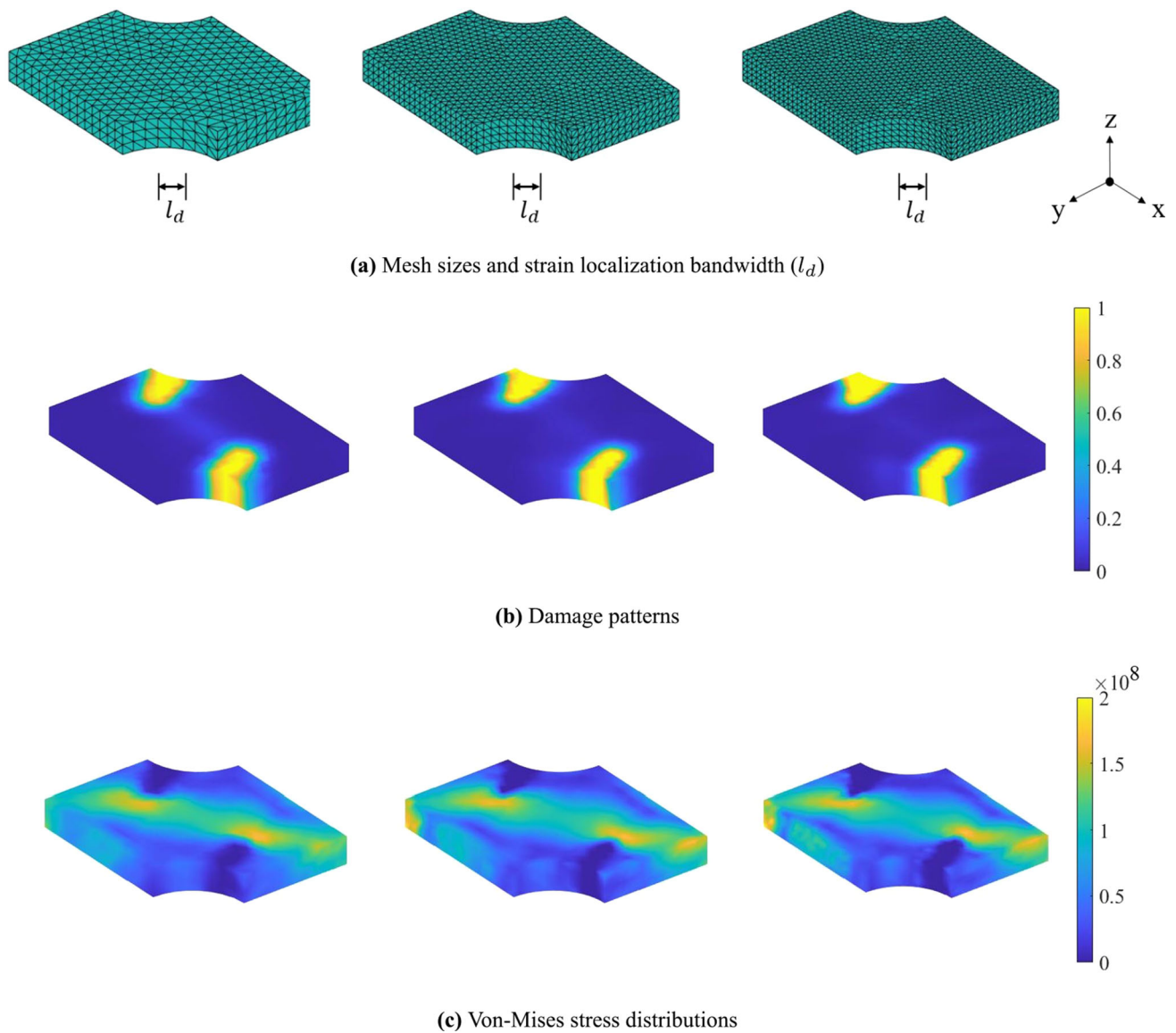


Fig. 18 Mesh convergence study of field variables by our FE-RNN: With the increase of the number of macro-elements in (a), both the damage patterns in (b) and stress distributions in (c) show convergence

Table 3 Computational time of the double notched model: The multiscale simulations of the double notched specimen are performed by the proposed FE-RNN model on a modest desktop with four CPU cores where the computational costs of one-time data generation and model training are reported in Table 2

Number of macro-elements	Multiscale simulation time
7,000	3.1×4 CPU-hour
18,000	7.5×4 CPU-hour
28,000	10.5×4 CPU-hour

deep learning model is an RNN trained on a database that is generated by efficiently sampling deformation path

spaces (spanned by strains) and obtaining the corresponding microstructural responses via a physics-based reduced-order model. To reduce the reliance on expensive data while increasing accuracy, we leverage two constraints while building our RNN. The first constraint is based on our energy analysis on thermodynamic consistency of microstructural deformation and we implement this constraint by adding a penalty term to our RNN’s loss function. The second constraint draws inspiration from the irreversible nature of damage processes and we implement it as a hard constraint by directly manipulating the temporal variation of the outputs within the RNN architecture. In addition, we incorporate the teacher forcing technique into our RNN model and demonstrate its impacts in both single and multiscale simulations.

We validate the accuracy of our model in both single and multiscale simulations that involve path-dependent deformations with hardening and softening. Importantly, we show that our model is accurate enough to provide reliable multiscale simulations with complex and cyclic loadings, particularly, we can reliably capture softening behaviors such that the solutions are post-failure convergent and mesh independent.

Our experiments reveal that while the costs of database generation and model training are considerable, these costs are amortized in online computations. For example, our data-driven model is about four orders of magnitude faster than the classic FE^2 approach in terms of CPU hours. Such high efficiency makes our model promising for many computationally intensive tasks that would require large computational resources (multi-core CPUs and GPUs) or need long simulation times.

The main objective of this work is to develop an RNN-based surrogate to accelerate multiscale damage analysis for metallic alloys with process-induced porosity. While the physics-based constraints that we use in developing this RNN improve its accuracy, the training dataset (generated with offline microstructural analyses) is the major driving factor that controls its behavior. That is, if a particular aspect of the physics is not present (explicitly or implicitly) in the dataset, our trained model cannot capture it.

Additionally, our surrogate is trained on a synthetic database which we believe is an important step before systematically incorporating experimental data (which are typically small) in a machine learning model. We plan to investigate the combination of synthetic and experimental data as they complement each other.

We also plan to extend our work in a number of different directions. First, minimization of inference error is critical especially for iterative solvers. While teacher-forcing is beneficial in single scale simulation, we are interested in leveraging this mechanism for online iterative multiscale computations where ground truth values are not available (simply adding noise during training does not seem to help based on our studies). Second, we plan to study the impacts of spatially varying material properties and microstructural morphology on the behaviors of macro components. However, adding such variations dramatically increases the dimension of sampling space and, in turn, the number of data points required for effective training. To reduce sampling costs in such scenarios, adaptive sampling strategies [61] for sequence learners or mechanistic neural networks that systematically couple data science with principles of mechanics [62, 63] need to be investigated. Lastly, we plan to make our model probabilistic for outer-loop applications such as uncertainty quantification [64] and material/structure design optimization [62, 65–68].

Acknowledgements The authors appreciate the supports from ACRC consortium, the feedback from the anonymous reviewers, and the helpful discussions with Dr. Ling Wu and Dr. Ludovic Noels. Ramin Bostanabad also acknowledges support from NSF (award number 2211908) and the Office of Naval Research (award number N00014-23-1-2485).

Declarations

Competing Interest The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Appendices

A Multiscale Continuum Damage Model

There are two popular approaches for modeling damage [69, 70]: (1) a discrete approach based on fracture mechanics which models displacement discontinuity at discontinuous interfaces, and (2) a continuous approach based on continuum mechanics which models damage as strain softening by inelastic strains. In our work, we adopt the latter approach to generate a database of microstructural responses. The continuous approach, however, often suffers from non-objectivity where localized damage bands become narrower and narrower upon mesh refinement, i.e., they are mesh dependent (this dependency is caused by ill-posed governing equations due to damage-induced non-positive definiteness). One way to address this issue is to modify the continuum constitutive law by introducing crack bandwidth via nonlocal functions.

However, integrating nonlocal functions within a multiscale damage model is quite difficult [18, 60]. This is because, on the one hand, a material characteristic length l_m is needed to stabilize ill-posed governing equations (caused by non-positive definite stiffness matrix) to address the mesh dependency at microscale. On the other hand, a macroscale characteristic length l_M is required to provide nonlocality for distant macro elements. Simultaneously imposing l_m and l_M helps to stabilize the multiscale damage model but is not physically realistic. If only l_m is imposed and l_M is neglected, the multiscale model merely transfers the damage-induced localization to the finer scale which is equivalent to a single scale model with l_m and contradicts with the purpose of multiscale modeling without properly transmitting microscale damage to macroscale [60].

One way to properly impose nonlocal functions in multiscale damage models is to introduce a material characteristic length in each RVE that accounts for the influence of damage on neighboring RVEs (and associated macro elements) [71]. We adopt another approach to properly impose both the micro and macro damage regularizations [26]: our micro damage model is regularized by a predefined material frac-

ture energy and its resulting effective damage parameter is subsequently regularized by a macro damage model via a nonlocal function for neighboring RVEs. We provide more details on our approach below.

We adopt an continuum damage model to simulate strain softening in ductile metals whose load-carrying capacity drops due to the degradation of yield stress and stiffness. To simulate the onset of softening, we choose ductile damage initiation criteria which models the effective strain at damage initiation, i.e., \bar{E}_d^{pl} , as a function of stress and strain states. We presume \bar{E}_d^{pl} is a constant and that damage begins when the equivalent plastic strain is equal or greater than it, i.e., $\bar{E}^{pl} \geq \bar{E}_d^{pl}$.

A major challenge of continuum damage models is the softening-induced non-positive definite stiffness matrix that results in slow solution convergence and negative wave speeds [60]. Specifically, the ill-posed problem causes equilibrium equations to lose objectivity with respect to mesh sizes by exhibiting spurious mesh sensitivity. In microstructural simulations, to address the lack of objectivity to mesh choices, we convert the stress–strain relation in the constitutive equation to a stress–displacement relation to formulate the micro-damage evolution after damage initiation as:

$$G_f = \int_{\bar{E}_0^{pl}}^{\bar{E}_f^{pl}} l_e S_y d\bar{E}^{pl} = \int_0^{\bar{u}_f^{pl}} S_y d\bar{u}^{pl} \tag{A-1}$$

where l_e indicates the element’s characteristic length in an arbitrary RVE, and G_f represents the dissipated energy that opens a unit area of crack after damage initiation. The equivalent plastic displacement \bar{u}^{pl} is the fracture work conjugate to the yield stress S_y from the onset of damage (with the effective plastic strain \bar{E}_0^{pl} and zero plastic displacement \bar{u}^{pl}) until the final failure (with the effective fracture strain \bar{E}_f^{pl} and the fracture displacement \bar{u}_f^{pl}). Using Equation (A-1) we define the microscale damage evolution based on an exponential form of the released energy [72] as:

$$D_m = 1 - \exp\left(-\frac{1}{G_f} \int_0^{\bar{u}_f^{pl}} S_y d\bar{u}^{pl}\right) \tag{A-2}$$

where D_m represents the microscale damage parameter that monotonically increases in the range of [0.0, 1.0]. We note that in our context of isotropic continuum damage, D_m is a scalar and it becomes a tensor in anisotropic damage models. In addition, we note that D_m approaches 1.0 asymptotically with infinitely large \bar{u}^{pl} in Equation (A-2). In practice, we set D_m as 1.0 when the dissipated energy exceeds $0.99G_f$. In our continuum damage model, we formulate a softening response of a micro point with an elasto-plastic behavior as:

$$S_m = (1 - D_m)S_m^0; \quad S_m^0 = \mathbb{C}^{el} : E_m^{el} = \mathbb{C}^{el} : (E_m - E_m^{pl}) \tag{A-3}$$

where S_m and S_m^0 are, respectively, the damaged stress and the reference stress that undergoes the same deformation path but in the absence of damage at a micro material point. \mathbb{C}^{el} represents the fourth-order elasticity tensor. E_m , E_m^{el} and E_m^{pl} are the microscale total strain, elastic strain, and plastic strain, respectively.

From the microscale stress, we can compute the effective stress via Equation (2). Additionally, we compute the RVE’s effective damage parameter [71] by:

$$D_M = 1 - \frac{\|S_M : S_M^0\|}{\|S_M^0 : S_M^0\|} \tag{A-4}$$

where the homogenized damage parameter D_M indicates the damage status of the RVE. Its value depends on the values of the effective stress S_M and the effective reference stress S_M^0 without damage; thus, it is clear that the effective damage parameter is not a function of homogenized plastic strains.

We proceed to constrain the effective damage parameter D_M via an integral-type non-local damage model to mitigate the spurious mesh dependency on the macroscale as:

$$\hat{D}_M(\mathbf{P}, \mathbf{P}') = \int_B \omega(\|\mathbf{P} - \mathbf{P}'\|) D_M(\mathbf{P}') d\mathbf{P}' \tag{A-5}$$

where $\hat{D}_M(\mathbf{P}, \mathbf{P}')$ is the non-local damage parameter at the macroscopic point \mathbf{P} surrounded by points \mathbf{P}' in the compact neighborhood B . $D_M(\mathbf{P}')$ represents the local damage parameter at \mathbf{P}' and ω indicates the non-local weighting function which depends on the distance $\|\mathbf{P} - \mathbf{P}'\|$ between the studied point and its supporting points. In this work, we define ω by a polynomial bell-shape function as:

$$\omega(\|\mathbf{P} - \mathbf{P}'\|) = \frac{\omega_\infty(\|\mathbf{P} - \mathbf{P}'\|)}{\int_B \omega_\infty(\|\mathbf{P} - \mathbf{P}'\|) d\mathbf{P}'} \tag{A-6a}$$

$$\omega_\infty(\|\mathbf{P} - \mathbf{P}'\|) = \left\langle 1 - \frac{4(\|\mathbf{P} - \mathbf{P}'\|)^2}{l_d^2} \right\rangle^2 \tag{A-6b}$$

where $\langle \dots \rangle$ is the Macauley bracket defined as $\langle x \rangle = \max(0, x)$, l_d denotes the strain localization bandwidth whose value represents the non-local interacting radius, and the support domain B is a sphere with a radius of $l_d/2$ in 3D models.

In our multiscale damage model, we compute the regularized macro stress, S , as:

$$S = (1 - \hat{D}_M)S_M^0 \tag{A-7}$$

where we assume the effective reference stress S_M^0 can be closely approximated by the effective damaged stress (S_M)

and the effective damage parameter (D_M) from microstructural analyses. That is:

$$\mathbf{S}_M^0 = \mathbf{S}_M / (1 - D_M) \quad (\text{A-8})$$

We can directly relate \mathbf{S} to the RVEs' effective damaged stress (\mathbf{S}_M) via $\mathbf{S} = \mathbf{S}_M(1 - \hat{D}_M)/(1 - D_M)$. These two stresses are identical if there are no macro nonlocal functions, i.e., $\hat{D}_M = D_M$. However, as discussed before, it is important to properly impose damage regularization on each scale to address the mesh dependency issue of continuum damage mechanics [18, 60]. In our model, we regularize micro damage by material fracture energy in Equation (A-2) and macro damage by nonlocal functions in Equation (A-5). Hence, the two stresses (\mathbf{S} and \mathbf{S}_M) in our case are directly related via a coefficient of $(1 - \hat{D}_M)/(1 - D_M)$.

We now demonstrate that the relation in Equation (A-8) is directly related to the definition of the effective damage parameter in Equation (A-4):

$$\mathbf{S}_M = (1 - D_M)\mathbf{S}_M^0 \quad (\text{A-9a})$$

$$\mathbf{S}_M : \mathbf{S}_M^0 = (1 - D_M)\mathbf{S}_M^0 : \mathbf{S}_M^0 \quad (\text{A-9b})$$

$$(\mathbf{S}_M : \mathbf{S}_M^0) / (\mathbf{S}_M^0 : \mathbf{S}_M^0) = (1 - D_M)(\mathbf{S}_M^0 : \mathbf{S}_M^0) / (\mathbf{S}_M^0 : \mathbf{S}_M^0) \quad (\text{A-9c})$$

$$(\mathbf{S}_M : \mathbf{S}_M^0) / (\mathbf{S}_M^0 : \mathbf{S}_M^0) = (1 - D_M) \quad (\text{A-9d})$$

where we can obtain the definition of the effective damage parameter as in Equation (A-4) since $0 \leq 1 - D_M \leq 1$.

B Hybrid Constitutive Integration

The non-positive definiteness of the stiffness matrix is the primary reason for the slow convergence of classic implicit time integration schemes that are used in continuum damage simulations. For illustration, consider the constitutive equation of an isotropic damage model integrated by an implicit backward-Euler integration scheme. Its algorithmic tangent operator at an arbitrary macroscopic IP can be written as:

$$\begin{aligned} \mathbb{C}_{n+1}^{alg} &= \frac{\partial \mathbf{S}_{n+1}}{\partial \mathbf{E}_{n+1}} = (1 - D_{n+1})\mathbb{C}^{el} \\ &\quad - \frac{S_{n+1} - H_n \bar{E}_{n+1}^{pl}}{(\bar{E}_{n+1}^{pl})^3} \mathbf{S}_{n+1}^0 \otimes \mathbf{S}_{n+1}^0 \end{aligned} \quad (\text{B-1})$$

where \mathbb{C}_{n+1}^{alg} , \bar{E}_{n+1}^{pl} , S_{n+1} , \mathbf{S}_{n+1}^0 and H_n represent the fourth-order algorithmic tangent operator, equivalent plastic strain, equivalent stress, referenced stress tensor, and softening modulus, respectively. The subscripts denote time steps and the symbol \otimes represents the cross product between tensors. Softening causes negative values for H_n which can

render \mathbb{C}_{n+1}^{alg} indefinite. A non-positive \mathbb{C}_{n+1}^{alg} leads to an ill-conditioned elemental stiffness matrix with near-zero or negative eigenvalues, and further deteriorates the global stiffness matrix in the element assembly process. Such ill-posed matrices dramatically reduce the efficiency of iterative solvers (e.g., Newton–Raphson methods) and often cause job abortion before final convergence.

To fundamentally resolve the convergence issue, we adopt a hybrid time integration scheme [26, 73] to integrate the governing equations of elasto-plastic and softening equations explicitly-implicitly. The basic idea of the hybrid integration is to maintain the positive definiteness of the system's algebraic tangent operator by separately integrating constitutive equations in two consecutive stages via explicit and implicit schemes. At the first stage, we explicitly extrapolate internal material state variables at time step $n + 1$ from step n to compute the explicit stress state $\tilde{\mathbf{S}}_{n+1}$ that balances the equilibrium equation between internal and external forces. At the second stage, we compute the implicit stress state \mathbf{S}_{n+1} based on the current strain state \mathbf{E}_{n+1} using the classic backward Euler method to update the trial stress tensor and yield functions for the next time step where the tangent operator between $\tilde{\mathbf{S}}_{n+1}$ and \mathbf{E}_{n+1} is kept positive definite.

For the elasto-plastic model, we choose the material state variable as the incremental plastic strain tensor $\Delta \tilde{\mathbf{E}}_{n+1}^{pl}$ such that $\tilde{\mathbf{S}}_{n+1}$ can be computed as:

$$\begin{aligned} \tilde{\mathbf{S}}_{n+1}(\Delta \tilde{\mathbf{E}}_{n+1}^{pl}) &= \tilde{\mathbf{S}}_{n+1}^{trial} - \mathbb{C}^{el} : \Delta \tilde{\mathbf{E}}_{n+1}^{pl} = \mathbb{C}^{el} : \mathbf{E}_{n+1} \\ &\quad - \mathbb{C}^{el} : \mathbf{E}_n^{pl} - \mathbb{C}^{el} : \Delta \tilde{\mathbf{E}}_{n+1}^{pl} \\ \Delta \tilde{\mathbf{E}}_{n+1}^{pl} &= \frac{\Delta t_{n+1}}{\Delta t_n} \Delta \mathbf{E}_n^{pl} \end{aligned} \quad (\text{B-2})$$

where \mathbf{E}_n^{pl} represents the implicit incremental plastic strain tensor at time step n , Δt_n and Δt_{n+1} indicate the lengths of time steps at two consecutive steps. The algorithmic tangent operator (under loading⁵) is therefore computed as:

$$\begin{aligned} \tilde{\mathbb{C}}_{n+1}^{alg} &= \frac{\partial \tilde{\mathbf{S}}_{n+1}(\Delta \tilde{\mathbf{E}}_{n+1}^{pl})}{\partial \mathbf{E}_{n+1}} \\ &= \frac{\partial (\mathbb{C}^{el} : \mathbf{E}_{n+1} - \mathbb{C}^{el} : \mathbf{E}_n^{pl} - \mathbb{C}^{el} : \Delta \tilde{\mathbf{E}}_{n+1}^{pl})}{\partial \mathbf{E}_{n+1}} = \mathbb{C}^{el} \end{aligned} \quad (\text{B-3})$$

In a similar manner, for isotropic continuum damage models, we choose the explicitly interpolated material state variable in the hybrid integration as the incremental plastic multiplier $\Delta \tilde{\lambda}_{n+1}$, i.e., $\Delta \tilde{\lambda}_{n+1} = (\Delta t_{n+1} / \Delta t_n) \Delta \lambda_n$. We can then write its explicit damaged stress and algorithmic tangent

⁵ $\tilde{\mathbb{C}}_{n+1}^{alg}$ is equal to the elastic modulus in unloading.

operator (under loading⁶) as:

$$\begin{aligned} \tilde{S}_{n+1} &= (1 - \tilde{D}_{n+1})S_{n+1}^0 = (1 - \tilde{D}_{n+1})\mathbb{C}^{el} : \mathbf{E}_{n+1}; \\ \tilde{D}_{n+1} &= \tilde{D}_{n+1}(D_n, \Delta\tilde{\lambda}_{n+1}) \end{aligned} \tag{B-4}$$

$$\tilde{\mathbb{C}}_{n+1}^{alg} = \frac{\partial \tilde{S}_{n+1}}{\partial \mathbf{E}_{n+1}} = (1 - \tilde{D}_{n+1})\mathbb{C}^{el} \tag{B-5}$$

where S_{n+1}^0 is the effective stress tensor, and \tilde{D}_{n+1} represents the explicit state of the damage variable which is a function of its previous implicit state D_n and the current explicit incremental plastic multiplier $\Delta\tilde{\lambda}_{n+1}$.

In the hybrid integration scheme, the loading tangent operators of the elasto-plastic model in Equation (B-3) and the damage model in Equation (B-5) are trivially equal to the elastic modulus \mathbb{C}^{el} and $(1 - \tilde{D}_{n+1})\mathbb{C}^{el}$. Hence, the hybrid integration scheme preserves the positive-definiteness of the governing equations and also allows to assemble the global stiffness matrix only once before online simulations. The global stiffness matrix remains constant for the elasto-plastic regime and only needs partial updates on matrix entries associated with the softening IPs by Equation (B-5). As softening is often highly localized in small regions, the global stiffness can be incrementally updated during the entire elasto-plastic-hardening-softening process [26]; saving significant memory footprints with robust convergence performance.

C Deflated Clustering Analysis

Simulation of microstructural softening via the classic FE² method involves demanding computational costs, which is prohibitive for generating big training data for machine learning models. To accelerate the database generation, we adopt our previously developed mechanistic ROM, i.e., deflated clustering analysis (DCA) [25, 26]. Its high efficiency comes from two facts: (1) the number of unknown variables in the system is dramatically reduced from a large number of finite elements to a few clusters by agglomerating elements via clustering as shown in Fig. 23, and (2) the algebraic equations of the reduced system contains much fewer close-to-zero eigenvalues that results in better convergence comparing to the classic FE system.

Our DCA utilizes k-means clustering, i.e., an unsupervised machine learning technique for data interpretation and grouping, to agglomerate neighboring elements into a set of interactive irregular-shape clusters. The clustering begins with feeding the coordinates of element centroids into a feature space where randomly scattered cluster seeds serve as initial cluster means. Clusters accept or reject elements by iteratively minimizing the within-cluster variance until all

⁶ $\tilde{\mathbb{C}}_{n+1}^{alg} = (1 - \tilde{D}_{n+1})\mathbb{C}^{el}$ in unloading.

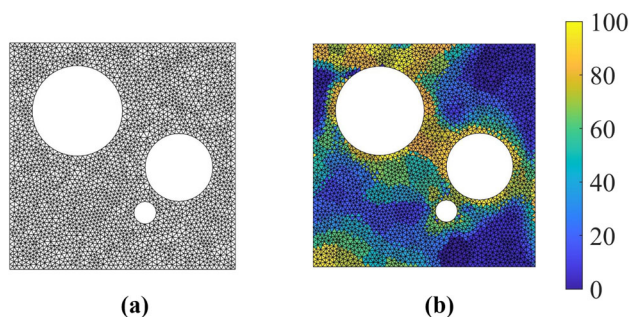


Fig. 19 Demonstration of clustering in ROM: The domain of a generic 2D RVE with 5, 000 elements in (a) are decomposed into 100 clusters in (b) where elements in the same cluster are assigned with the same color

elements are assigned to a cluster. The clustering procedure can be mathematically stated as a minimization problem as:

$$C = \min_C \sum_{l=1}^k \sum_{n \in C^l} \|\varphi_n - \bar{\varphi}_l\|^2 \tag{C-1}$$

where C represents the k clusters with $C = \{C^1, C^2, \dots, C^k\}$, φ_n and $\bar{\varphi}_l$ indicate the coordinates of the centroid of the n^{th} element and the mean of the coordinates of the l^{th} cluster, respectively. A clustering example is illustrated in Fig. 23 where the discrete domain of a 2D generic RVE with 5, 000 elements is decomposed into 100 clusters.

We construct clustering-based reduced mesh via Delaunay triangularization by connecting cluster centroids where the topological relations between clusters are preserved from the original FE mesh. By assuming the motions of cluster centroids are directly related to clustering nodes, we can compute the nodal displacements via polynomial augmented radian point interpolation [74] as:

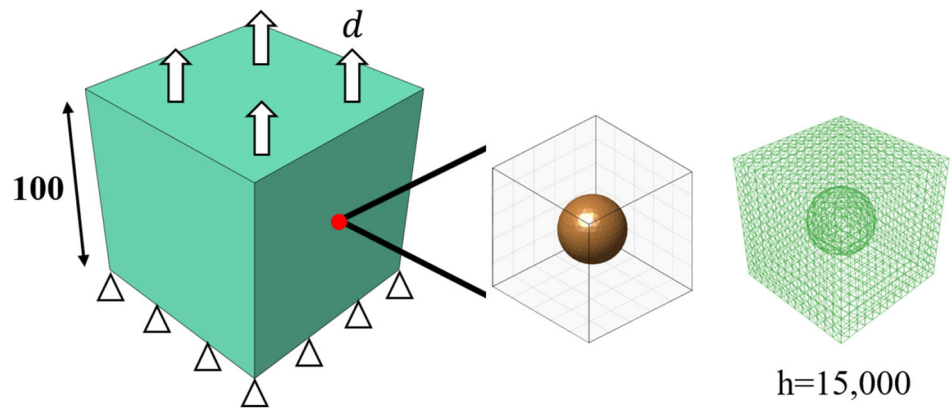
$$u_c = \mathbf{R}\mathbf{a} + \mathbf{Z}\mathbf{b} \tag{C-2}$$

where u_c represents the displacements of cluster centroids. \mathbf{a} is the coefficient vector of the radial basis function matrix \mathbf{R} , and \mathbf{b} is the coefficient vector of the polynomial basis matrix \mathbf{Z} . Meanwhile, the radial coefficient and the polynomial basis need to satisfy the following equation for every node per cluster and every polynomial basis function to ensure solution uniqueness [74] as:

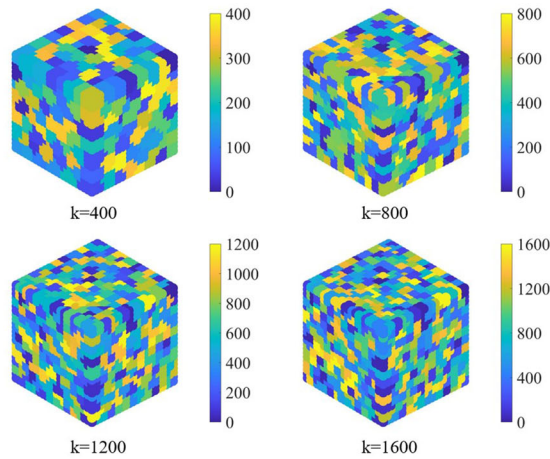
$$\mathbf{Z}\mathbf{a} = \mathbf{0} \tag{C-3}$$

The displacements of cluster centroids are augmented with rotational degrees of freedom to represent the six rigid body motions in a 3D deflation space [75], including three translations and three rotations. Upon the completion of a non-linear analysis on the reduced mesh, the displacement solutions can

Fig. 20 Multiscale cube model: **a** Every integration point of the macro-cube model is associated with a porous RVE; and **b** The RVE domain is discretized by different numbers of clusters



(a) Geometry, dimensions (unit: mm), and boundary conditions of the macro-model; Geometry and the finite element mesh of the porous RVE



(b) RVE clustering

be projected back to the original FE mesh by:

$$\mathbf{u}_i^j = \mathbf{W}_i^j \boldsymbol{\lambda}_j \tag{C-4}$$

where \mathbf{u}_i^j represents the displacement vector at the i^{th} node in the j^{th} cluster. In addition, $\boldsymbol{\lambda}_j$ is the rigid body motion of the centroid of the j^{th} cluster, while the \mathbf{W}_i^j indicates the deflation matrix for the i^{th} node in the j^{th} cluster as:

$$\boldsymbol{\lambda}_j = [u_{jx}, u_{jy}, u_{jz}, \theta_{jx}, \theta_{jy}, \theta_{jz}]^T;$$

$$\mathbf{W}_i^j = \begin{bmatrix} 1 & 0 & 0 & z_i^j & -y_i^j \\ 0 & 1 & 0 & -z_i^j & x_i^j \\ 0 & 0 & 1 & y_i^j & -x_i^j & 0 \end{bmatrix} \tag{C-5}$$

where u_{jx} and θ_{jx} are the displacement and rotation of the j^{th} cluster along x axis, and the (x_i^j, y_i^j, z_i^j) are the relative 3D coordinates of the i^{th} node with respect to the centroid of

the j^{th} cluster. By assuming all elements in the same cluster share identical stress and strain fields, microstructural effective responses can be reproduced in a highly efficient manner such that the unknown variables are dramatically decreased from FE system that accounts for distinct field variables per element to the reduced system with much fewer distinct solutions per cluster.

To demonstrate the efficacy of our DCA, we compare its simulation results on a 3D multiscale cube against the classic FE² method in Fig. 20. The macro-cube is fully constrained at its bottom surface, and it is subject to an upward extension on the top surface with $d = 7$ mm. The cube is meshed with 12 tetrahedral elements of reduced-integration (one IP at the center of each tetrahedron). We assume each macro-IP is associated with the same porous RVE containing one spherical pore in the middle as shown in Fig. 20a.

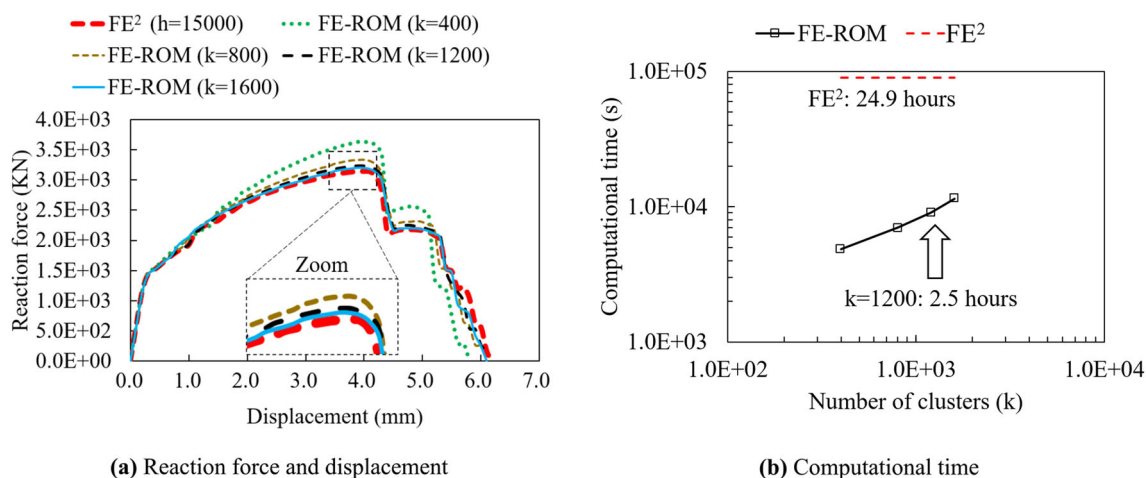


Fig. 21 Results of the multiscale cube model: **a** Comparison of the softening load–displacement curves between FE^2 and FE-ROM with different clusters; and **b** Comparison of computational time

To determine the number of clusters for a given problem (for any clustering-based ROM, e.g., DCA, SCA, or SCA’s variants), we can perform a quick preliminary convergence study where we gradually increase the number of clusters and determine the minimum number of clusters above which the results insignificantly change. This convergence study can also be done by comparing the results of the ROM to that of direction numerical solutions (DNS). Yet another method is to formulate a data-driven inverse optimization problem [59] where the cluster number is considered as an optimization variable. In this work, we carry out a convergence study to ensure our ROM’s solutions do not change as the number of clusters increase and that they are consistent with the DNS, i.e., FE^2 . Specifically, we apply four clustering levels (k) of 400, 800, 1, 200 and 1, 600 to an RVE meshed with 15, 000 elements and investigate the effects of k on the RVE’s effective softening behaviors, see Fig. 20.

We compare the reaction force-displacement curves from FE^2 and FE-ROM in Fig. 21a. By considering the FE^2 solutions as the benchmark, we observe that: (1) the FE-ROM solutions with $k = 400$ slightly overestimate the component’s strength as insufficient clustering in the RVE artificially strengthens the material [23, 25]; and (2) as k increases, the FE-ROM responses (especially the post-failure behaviors) become closer and closer to the benchmark. Specifically, we observe that when k increases to 1, 200 and 1, 600, FE-ROMs achieve sufficiently accurate results compared to FE^2 .

We compare the computational costs of the different solvers in Fig. 21b. While all experiments are performed on an HPC by paralleling 60 CPU cores with 360 GB RAM, the clock time of FE^2 is the longest (about 24.9 hours). The clock time of the ROM with 1, 200 and 1, 600 clusters is about 2.5 and 3.2 hours, resulting in the acceleration fac-

tors of 9.9 and 7.8, respectively. Considering the fact that the ROM with $k = 1, 200$ is about 28% faster than its counterpart with $k = 1, 600$ while achieving similar accuracy, we adopt $k = 1, 200$ while building the training dataset in Sect. 4.

For efficient generation of (micro)structure-performance datasets, we note that many other ROMs can also be used for porous microstructural analyses. For example, self-consistent analysis (SCA) [23, 76, 77] and virtual clustering analysis (VCA) [24] can achieve highly efficient and accurate microstructural homogenization results by treating pores as a soft material with the 0.1% modulus of matrix materials [78]. Another method is the FEM-cluster-based analysis (FCA) [79] where the Hill-Mandel theorem is replaced with the energy equivalence theorem without filling pores with reference material properties. As our focus in this paper is on building the deep learning model that can faithfully surrogate microstructural analyses, we use our in-house DCA package and plan to leverage other methods such as SCA in our future works.

D Gated Recurrent Unit

To alleviate vanishing and exploding gradient issues of RNNs in processing long sequential data, long short term memory (LSTM) and gated recurrent unit (GRU) are typically used. GRU is a variant of the LSTM that, while providing similar accuracy, is more parsimonious and hence computationally more efficient. It is for this reason that we choose GRU as the memory cell in our proposed RNN architecture as in Fig. 4.

To demonstrate the working mechanism of GRUs, we three interconnected cells of a GRU layer in Fig. 22. In a GRU layer, a typical cell at an arbitrary time step t generates predictions \hat{y}_t and internal memory-like hidden variables h_t

Fig. 22 Architectures of GRU layer and cells: The internal structure and mathematical operations are demonstrated in the GRU cell at time step t

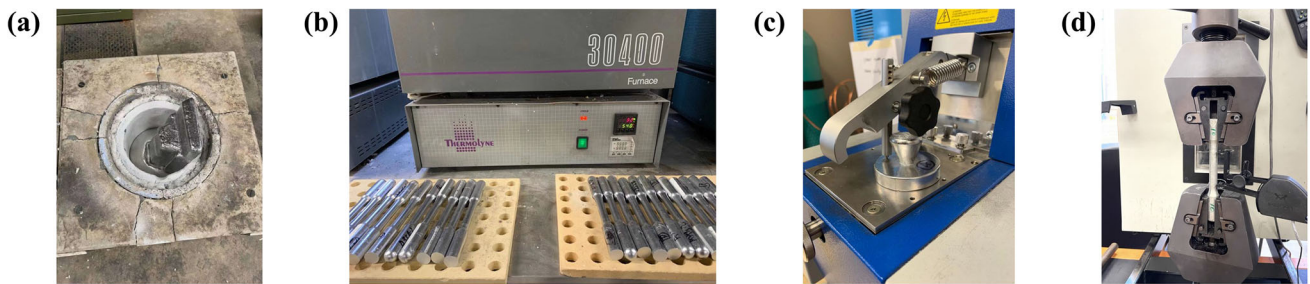
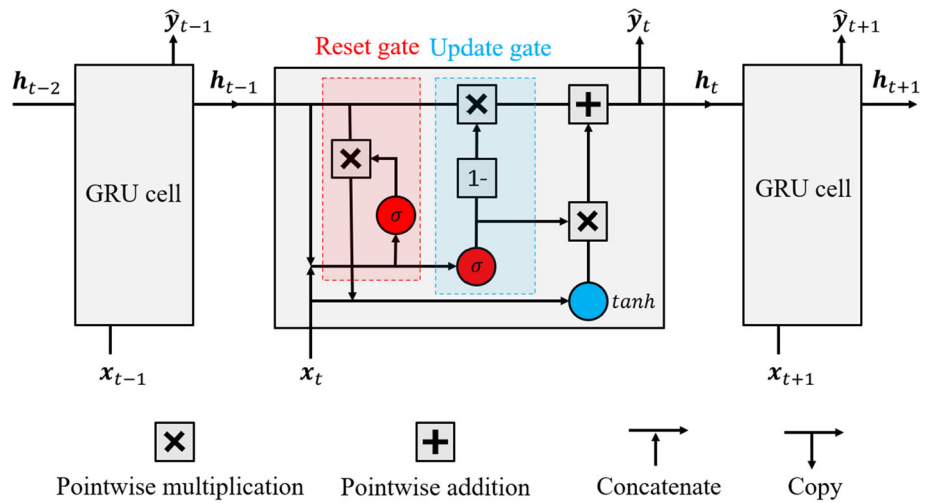


Fig. 23 Experimental characterization of our aluminum alloy A356: **a** A356 ingots are melted in a high-temperature furnace with degassing to remove porosity; **b** Heat treatment of cast tensile bars; **c** Composition analysis; and **d** Tensile tests of the cast alloys

after reading in the current inputs x_t and the hidden variables h_{t-1} from the previous cell. Compared to the RNN cell in Fig. 3b, the GRU cell uses reset and update gates to regulate its internal information flow. The reset gate r_t reads x_t and h_{t-1} to determine the candidate hidden state \tilde{h}_t by filtering out less important information passing from the previous cell. Its operations include:

$$r_t = \sigma(W_{hr}h_{t-1} + W_{xr}x_t + b_r) \tag{D-1a}$$

$$\tilde{h}_t = \tanh(r_t \odot W_{h\tilde{h}}h_{t-1} + W_{x\tilde{h}}x_t + b_{\tilde{h}}) \tag{D-1b}$$

where σ is the sigmoid activation function that returns a value in the range of $[0, 1]$, \tanh is the hyperbolic tangent function, and \odot represents the Hadamard product. W_{hr} , W_{xr} , $W_{h\tilde{h}}$, $W_{x\tilde{h}}$ are the weight matrices associated with the hidden state, the input state, the hidden-to-candidate hidden state and the input-to-candidate hidden state, respectively. b_r and $b_{\tilde{h}}$ are the biases applied to the sigmoid function in the reset gate and the hyperbolic tangent function, respectively.

The update gate (which has its weights and biases) similarly operates on x_t and h_{t-1} : it linearly interpolates the previous hidden state h_{t-1} and the candidate hidden state \tilde{h}_t to update the memory-like hidden state h_t which is then

passed to the next cell:

$$u_t = \sigma(W_{hu}h_{t-1} + W_{xu}x_t + b_u) \tag{D-2a}$$

$$h_t = u_t \odot h_{t-1} + (1 - u_t) \odot \tilde{h}_t + b_h \tag{D-2b}$$

where W_{hu} and W_{xu} are the weights applied onto the hidden state and input state in the update gate. b_u and b_h are the two biases associated to the sigmoid function and the generation of the current hidden state. The cell output at the current time step \hat{y}_t is then obtained by linearly transforming the hidden state:

$$\hat{y}_t = W_{hy}h_t + b_y \tag{D-3}$$

where W_{hy} and b_y are the weights and biases associated with the current output state \hat{y}_t . We note that all the weights and biases of the GRU networks are iteratively updated by BPTT during training.

E Experimental Material Characterization

For the microstructural simulations in Appendix C we assume the microstructure only contains porosity and the

matrix material (i.e., aluminum alloy A356). So, in this section, we briefly discuss the experimental characterization process that can be used to obtain the effective elastoplastic and damage properties of the matrix material, see Fig. 23. Our experiment consists of several steps. In the first step, we melt aluminum A356 ingots in a furnace which is preheated to about 800° C. During the melting process, we apply degassing [80] to remove gases (e.g., hydrogen contents) and gas-induced porosity before casting as tensile coupons. In the second step, we apply a standard T6 heat treatment to improve the A356 alloy's strength and toughness. The heat treatment involves a high temperature treatment at 540° C for 8 hours to dissolve alloy elements into aluminum matrix, a quenching process to freeze alloy elements within the solid solution, and an artificial aging process at about 155° C for 3.5 hours to precipitate alloy elements and form grain structures. We also perform composition analysis and find that our A356 alloy contains about 92.05% aluminum (weight fraction), 6.72% silicon, 0.09% steel, 0.0028% magnesium, and other alloy elements. In the third step, we use X-ray computed tomography (CT) to inspect the porosity defect in tensile coupons to ensure the cast alloy is free of pores. Finally, we perform the tensile test on the tensile coupons and measure their averaged elastoplastic and damage parameters (which are provided in Sect. 4.1).

References

1. Feyel Frédéric, Chaboche Jean-Louis (2000) FE2 multiscale approach for modelling the elastoviscoplastic behaviour of long fibre SiC/Ti composite materials. *Comput Methods Appl Mech Eng* 183.3–4:309–330
2. Pascale Kanouté DP, Chaboche Boso Jean-Louis, Schrefler BA (2009) Multiscale methods for composites: a review. *Archiv Comput Methods Eng* 16(1):31–75
3. Jian-Ying Wu, Nguyen Vinh Phu, Nguyen Chi Thanh, Sutula Danas, Sinaie Sina, Bordas Stéphane PA (2020) Phase-field modeling of fracture. *Adv Appl Mech* 53:1–183
4. Griffith Alan Arnold (1921) VI. The phenomena of rupture and flow in solids. *Philos Trans Royal Soc London Ser A Contain Papers Math Phys Character* 221:163–198
5. Dugdale Donald S (1960) Yielding of steel sheets containing slits. *J Mech Phys Solids* 8(2):100–104
6. Bouchard Pierre-Olivier, Bay François, Chastel Yvan (2003) Numerical modelling of crack propagation: automatic remeshing and comparison of different criteria. *Comput Methods Appl Mech Eng* 192.35–36:3887–3908
7. Vinh Phu Nguyen and Hung Nguyen-Xuan (2013) High-order B-splines based finite elements for delamination analysis of laminated composites. *Compos Struct* 102:261–275
8. Jian-Ying Wu (2011) Unified analysis of enriched finite elements for modeling cohesive cracks. *Comput Methods Appl Mech Eng* 200.45–46:3031–3050
9. Moës Nicolas, Dolbow John, Belytschko Ted (1999) A finite element method for crack growth without remeshing. *Int J Numer Methods Eng* 46.1:131–150
10. Moës Nicolas, Gravouil Anthony, Belytschko Ted (2002) Non-planar 3D crack growth by the extended finite element and level sets—Part I: mechanical model. *Int J Numer Methods Eng* 53.11:2549–2568
11. Rashid Yrn R (1968) Ultimate strength analysis of prestressed concrete pressure vessels. *Nuclear Eng Design* 7.4:334–344
12. Cervera Miguel, Jian-Ying Wu (2015) On the conformity of strong, regularized, embedded and smeared discontinuity approaches for the modeling of localized failure in solids. *Int J Solids Struct* 71:19–38
13. Krajcinovic Dusan (1989) Damage mechanics. *Mech Mater* 8.2–3:117–197
14. Jirásek Milan (2007) Mathematical analysis of strain localization. *Revue européenne de génie civil* 11.7–8:977–991
15. Simo Juan C, Ju JW (1987) Strain-and stress-based continuum damage models-I. Formulation. *Int J Solids Struct* 23(7):821–840
16. De Borst R, Sluys LJ (1991) Localisation in a Cosserat continuum under static and dynamic loading conditions. *Comput Methods Appl Mech Eng* 90.1–3:805–827
17. Bazant Zdenek P, Belytschko Ted B, Chang Ta-Peng et al (1984) Continuum theory for strain-softening. *J Eng Mech* 110.12:1666–1692
18. Bazant Zdenek P, Jirásek Milan (2002) Nonlocal integral formulations of plasticity and damage: survey of progress. *J Eng Mech* 128.11:1119–1149
19. Poh Leong Hien, Sun Gang (2017) Localizing gradient damage model with decreasing interactions. *Int J Numer Methods Eng* 110.6:503–522
20. Bram Vandoren, Simone A (2018) Modeling and simulation of quasi-brittle failure with continuous anisotropic stress-based gradient-enhanced damage models. *Comput Methods Appl Mech Eng* 332:644–685
21. Dvorak George J (1992) Transformation field analysis of inelastic composite materials. *Proc Royal Soc London Ser A Math Phys Sci* 437(1900):311–327
22. Roussette Sophie, Michel Jean-Claude, Suquet Pierre (2009) Nonuniform transformation field analysis of elastic-viscoplastic composites. *Compos Sci Technol* 69.1:22–27
23. Liu Zeliang, Bessa MA, Liu Wing Kam (2016) Self-consistent clustering analysis: an efficient multi-scale scheme for inelastic heterogeneous materials. *Comput Methods Appl Mech Eng* 306:319–341
24. Tang Shaoqiang, Zhang Lei, Liu Wing Kam (2018) From virtual clustering analysis to self-consistent clustering analysis: a mathematical study. *Comput Mech* 62.6:1443–1460
25. Deng Shiguang, Soderhjelm Carl, Apelian Diran, Bostanabad Ramin (2022) Reduced-order multiscale modeling of plastic deformations in 3D alloys with spatially varying porosity by deflated clustering analysis. *Computat Mech* 70.3:517–548
26. Shiguang Deng, Diran Apelian, Ramin Bostanabad (2023) Adaptive spatiotemporal dimension reduction in concurrent multiscale damage analysis. *Computat Mech* 72:1–33
27. Planas R, Oune N, Bostanabad R (2021) Evolutionary Gaussian processes. *J Mech Design* 143(11):111703. <https://doi.org/10.1115/1.4050746>
28. Oune N, Bostanabad R (2021) Latent map Gaussian processes for mixed variable metamodeling. *Comput Methods Appl Mech Eng* 387:114128. <https://doi.org/10.1016/j.cma.2021.114128>
29. Chen W, Iyer A, Bostanabad R (2022) Data centric design: a new approach to design of microstructural material systems. *Engineering* 10:89–98. <https://doi.org/10.1016/j.eng.2021.05.022>
30. Zanjani Foumani Zahra, Mehdi Shishehbor, Amin Yousefpour, Ramin Bostanabad (2023) Multi-fidelity Costaware Bayesian optimization. *Comput Methods Appl Mech Eng* 407:115937. <https://doi.org/10.1016/j.cma.2023.115937>
31. Loujaine Mehrez, Jacob Fish, Venkat Aitharaju, Rodgers Will R, Roger Ghanem (2017) A PCE-based multiscale framework for the characterization of uncertainties in complex systems. *Com-*

- put Mech 61(1–2):219–236. <https://doi.org/10.1007/s00466-017-1502-4>. (ISSN: 0178-7675 1432-0924)
32. Carlos Mora, Tammer Eweis-Labolle Jonathan, Tyler Johnson, Likith Gadde, Ramin Bostanabad (2023) Probabilistic neural data fusion for learning from an arbitrary number of multi-fidelity data sets. *Comput Methods Appl Mech Eng* 415:116207. <https://doi.org/10.1016/j.cma.2023.116207>
 33. Jones RE, Templeton JA, Sanders CM, Ostien JT (2018) Machine learning models of plastic flow based on representation theory. *Comput Model Eng Sci* 117:309–342. <https://doi.org/10.31614/cmescs.2018.04285>
 34. Furukawa Tomonari, Yagawa Genki (1998) Implicit constitutive modelling for viscoplasticity using neural networks. *Int J Numer Methods Eng* 43.2:195–219
 35. Furukawa Tomonari, Hoffman Mark (2004) Accurate cyclic plastic analysis using a neural network material model. *Eng Anal Bound Elem* 28.3:195–204
 36. Fernández Mauricio, Rezaei Shahed, Mianroodi Jaber Rezaei, Fritzen Felix, Reese Stefanie (2020) Application of artificial neural networks for the prediction of interface mechanics: a study on grain boundary constitutive behavior. *Adv Model Simul Eng Sci* 7.1:1–27
 37. Xiaoxin Lu, Yvonne Julien, Detrez Fabrice, Bai Jinbo (2017) Multiscale modeling of nonlinear electric conductivity in graphene-reinforced nanocomposites taking into account tunnelling effect. *J Comput Phys* 337:116–131
 38. Mianroodi Jaber Rezaei, Siboni Nima H, Raabe Dierk (2021) Teaching solid mechanics to artificial intelligence—A fast solver for heterogeneous materials. *NPJ Comput Mater* 7.1:1–10
 39. Haghighat Ehsan, Raissi Maziar, Moure Adrian, Gomez Hector, Juanes Ruben (2021) A physics-informed deep learning framework for inversion and surrogate modeling in solid mechanics. *Comput Methods Appl Mech Eng* 379:113741
 40. Peivaste Iman, Siboni Nima H, Alahyarizadeh Ghasem, Ghaderi Reza, Svendsen Bob, Raabe Dierk, Mianroodi Jaber Rezaei (2022) Machine-learning-based surrogate modeling of microstructure evolution using Phasefield. *Comput Mater Sci* 214:111750
 41. Mozaffar M, Bostanabad R, Chen W, Ehmann K, Cao Jian, Bessa MA (2019) Deep learning predicts pathdependent plasticity. *Proc Natl Acad Sci* 116.52:26414–26420
 42. Wang Kun, Sun WaiChing (2018) A multiscale multi-permeability poroplasticity model linked by recursive homogenizations and deep learning. *Comput Methods Appl Mech Eng* 334:337–380
 43. Ling Wu, Kilingar Nanda Gopala, Noels Ludovic et al (2020) A recurrent neural network-accelerated multi-scale model for elastoplastic heterogeneous materials subjected to random cyclic and non-proportional loading paths. *Comput Methods Appl Mech Eng* 369:113234
 44. Ghavamian F, Simone A (2019) Accelerating multiscale finite element simulations of history-dependent materials using a recurrent neural network. *Comput Methods Appl Mech Eng* 357:112594
 45. Logarzo Hernan J, Capuano German, Rimoli Julian J (2021) Smart constitutive laws: inelastic homogenization through machine learning. *Comput Methods Appl Mech Eng* 373:113482
 46. Otero Fermin, Oller Sergio, Martínez Xavier (2018) Multiscale computational homogenization: review and proposal of a new enhanced-first-order method. *Archiv Comput Methods Eng* 25(2):479–505
 47. Tang Shaoqiang, Yang Yang (2021) Why neural networks apply to scientific computing? *Theor Appl Mech Lett* 11(3):100242
 48. Hornik Kurt, Stinchcombe Maxwell, White Halbert (1989) Multilayer feedforward networks are universal approximators. *Neural Netw* 2(5):359–366
 49. Lipton Zachary C, Berkowitz John, Elkan Charles (2015) A critical review of recurrent neural networks for sequence learning. In: arXiv preprint [arXiv:1506.00019](https://arxiv.org/abs/1506.00019)
 50. Hanin Boris (2018) Which neural net architectures give rise to exploding and vanishing gradients?. In: *Advances in neural information processing systems* vol 31
 51. Staudemeyer Ralf C, Morris Eric Rothstein (2019) Understanding LSTM—A tutorial into long short-term memory recurrent neural networks. In: arXiv preprint [arXiv:1909.09586](https://arxiv.org/abs/1909.09586)
 52. Karpathy Andrej, Johnson Justin, Fei-Fei Li (2015) Visualizing and understanding recurrent networks. In: arXiv preprint [arXiv:1506.02078](https://arxiv.org/abs/1506.02078)
 53. Silhavy Miroslav (2013) *The mechanics and thermodynamics of continuous media*. Springer, Berlin
 54. Yang Han, Sinha Sumeet Kumar, Feng Yuan, McCallen David B, Jeremić Boris (2018) Energy dissipation analysis of elastic-plastic materials. *Comput Methods Appl Mech Eng* 331:309–326
 55. Feigenbaum Heidi P, Dafalias Yannis F (2007) Directional distortional hardening in metal plasticity within thermodynamics. *Int J Solids Struct* 44.22–23:7526–7542
 56. Xiang Zixue, Peng Wei, Liu Xu, Yao Wen (2022) Self-adaptive loss balanced physics-informed neural networks. *Neurocomputing* 496:11–34
 57. Márquez-Neila Pablo, Salzmann Mathieu, Fua Pascal (2017) Imposing hard constraints on deep networks: Promises and limitations. In: arXiv preprint [arXiv:1706.02025](https://arxiv.org/abs/1706.02025)
 58. Goodfellow Ian, Bengio Yoshua, Courville Aaron (2016) *Deep learning*. MIT press, Cambridge
 59. Deng Shiguang, Mora Carlos, Apelian Diran, Bostanabad Ramin (2022) Data-driven calibration of Multifidelity multiscale fracture models via latent map Gaussian Process. *J Mech Design* 145(1):011705
 60. Bazant Zdenek P (2010) Can multiscale-multiphysics methods predict softening damage and structural failure? *Int J Multiscale Comput Eng* 8(1):61–67
 61. Bengio Samy, Vinyals Oriol, Jaitly Navdeep, Shazeer Noam (2015) Scheduled sampling for sequence prediction with recurrent neural networks. In: *Advances in neural information processing systems* vol 28
 62. Li Hengyang, Kafka Orion L, Gao Jiaying, Cheng Yu, Nie Yinghao, Zhang Lei, Tajdari Mahsa, Shan Tang Xu, Guo Gang Li et al (2019) Clustering discretization methods for generation of material performance databases in machine learning and design optimization. *Comput Mech* 64:281–305
 63. Liu Daoping, Hang Yang KI, Elkhodary Shan Tang, Liu Wing Kam, Guo Xu (2022) Mechanistically informed data-driven modeling of cyclic plasticity via artificial neural networks. *Comput Methods Appl Mech Eng* 393:114766
 64. Bostanabad Ramin, Liang Biao, Gao Jiaying, Liu Wing Kam, Cao Jian, Zeng Danielle, Xuming Su, Hongyi Xu, Li Yang, Chen Wei (2018) Uncertainty quantification in multiscale simulation of woven fiber composites. *Comput Methods Appl Mech Eng* 338:506–532
 65. Osanov Mikhail, Guest James K (2016) Topology optimization for architected materials design. *Annu Rev Mater Res* 46:211–233
 66. Zheng-Dong Ma, Noboru Kikuchi, Christophe Pierre, Basavaraju R (2006) Multidomain topology optimization for structural and material designs. *J. Appl. Mech.* 73(4):565–573
 67. Deng Shiguang, Suresh Krishnan (2016) Multi-constrained 3D topology optimization via augmented topological level-set. *Comput Struct* 170:1–12
 68. Deng Shiguang, Suresh Krishnan (2015) Multi-constrained topology optimization via the topological sensitivity. *Struct Multidiscip Optim* 51(5):987–1001
 69. Oliver Javier (1989) A consistent characteristic length for smeared cracking models. *Int J Numer Methods Eng* 28(2):461–474
 70. Oliver Javier, Huespe Alfredo Edmundo, Pulido MDG, Chaves E (2002) From continuum mechanics to fracture mechanics: the strong discontinuity approach. *Eng Fract Mech* 69.2:113–136

71. Liu Zeliang, Fleming Mark, Liu Wing Kam (2018) Microstructural material database for self-consistent clustering analysis of elastoplastic strain softening materials. *Comput Methods Appl Mech Eng* 330:547–577
72. Smith Michael (2009) ABAQUS standard user's manual. In: Dassault Systèmes Simulia Corp, Version 6.9
73. Oliver Javier, Huespe Alfredo Edmundo, Cante JC (2008) An implicit/explicit integration scheme to increase computability of non-linear material and contact/friction problems. *Comput Methods Appl Mech Eng* 19.721–24:1865–1889
74. Liu Gui-Rong (2009) Meshfree methods: moving beyond the finite element method. CRC Press, Boca Raton
75. Jönsthövel TB, Van Gijzen MB, Vuik C, Kasbergen C, Scarpas A (2009) Preconditioned conjugate gradient method enhanced by deflation of rigid body modes applied to composite materials. *Comput Model Eng Sci (CMES)* 47.2:97
76. Saha Sourav, Kafka Orion L, Ye Lu, Cheng Yu, Liu Wing Kam (2021) Macroscale property prediction for additively manufactured in625 from microstructure through advanced homogenization. *Integr Mater Manuf Innov* 10:360–372
77. Kafka Orion L, Cheng Yu, Cheng Puikui, Wolff Sarah J, Bennett Jennifer L, Garboczi Edward J, Cao Jian, Xiao Xianghui, Liu Wing Kam (2022) X-ray computed tomography analysis of pore deformation in IN718 made with directed energy deposition via in-situ tensile testing. *Int J Solids Struct* 256:111943
78. Yang Yang, Zhang Lei, Tang Shaoqiang (2022) A comparative study of cluster-based methods at finite strain. *Acta Mechanica Sinica* 38(4):421153
79. Nie Yinghao, Li Zheng, Cheng Gengdong (2021) Efficient prediction of the effective nonlinear properties of porous material by FEM-Cluster based Analysis (FCA). *Comput Methods Appl Mech Eng* 383:113921
80. Dispinar D, Akhtar Shahid, Nordmark Arne, Di Sabatino Marisa, Arnberg LJMS (2010) Degassing, hydrogen and porosity phenomena in A356. *Mater Sci Eng A* 527.16–17:3719–3725

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.