



HAL
open science

Deep reinforcement learning for the control of conjugate heat transfer

Elie Hachem, Hassan Ghraieb, Jonathan Viquerat, Aurélien Larcher, P. Meliga

► **To cite this version:**

Elie Hachem, Hassan Ghraieb, Jonathan Viquerat, Aurélien Larcher, P. Meliga. Deep reinforcement learning for the control of conjugate heat transfer. *Journal of Computational Physics*, 2021, 436, pp.110317. 10.1016/j.jcp.2021.110317 . hal-03027923v2

HAL Id: hal-03027923

<https://hal.science/hal-03027923v2>

Submitted on 17 Nov 2021

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

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

Deep reinforcement learning for the control of conjugate heat transfer

E. Hachem^{a,*}, H. Ghraieb^a, J. Viquerat^a, A. Larcher^a, P. Meliga^a

^a*MINES ParisTech, PSL Research University, Centre de mise en forme des matériaux (CEMEF), CNRS UMR 7635, 06904 Sophia Antipolis Cedex, France*

Abstract

This research gauges the ability of deep reinforcement learning (DRL) techniques to assist the control of conjugate heat transfer systems governed by the coupled Navier–Stokes and heat equations. It uses a novel, “degenerate” version of the proximal policy optimization (PPO) algorithm, intended for situations where the optimal policy to be learnt by a neural network does not depend on state, as is notably the case in optimization and open-loop control problems. The numerical reward fed to the neural network is computed with an in-house stabilized finite elements environment combining variational multi-scale (VMS) modeling of the governing equations, immerse volume method, and multi-component anisotropic mesh adaptation. Several test cases of natural and forced convection in two and three dimensions are used as testbed for developing the methodology. The approach successfully alleviates the natural convection induced enhancement of heat transfer in a two-dimensional, differentially heated square cavity controlled by piece-wise constant fluctuations of the sidewall temperature. It also proves capable of improving the homogeneity of temperature across the surface of two and three-dimensional hot workpieces under impingement cooling. Various cases are tackled, in which the position of multiple cold air injectors is optimized relative to a fixed workpiece position. The flexibility of the numerical framework makes it tractable to solve also the inverse problem, i.e., to optimize the workpiece position relative to a fixed injector distribution. The obtained results showcase the potential of the method for black-box optimization of practically meaningful computational fluid dynamics (CFD) conjugate heat transfer systems. More significantly, they stress how DRL can reveal unanticipated solutions or parameter relations (as the optimal workpiece position under symmetrical actuation turns to be offset from the symmetry axis), in addition to being a tool for optimizing searches in large parameter spaces.

Keywords: Deep Reinforcement Learning; Artificial Neural Networks; Conjugate heat transfer; Computational fluid dynamics; Thermal control.

1. Introduction

Thermal control, defined as the ability to finesse the thermal properties of a volume of fluid (and of the solid objects inside) into a certain desired state, is a field of tremendous societal and economical importance. For instance, heat/cool exchangers are used in a broad range of industrial applications to regulate process temperatures by heat or cool transfer between fluid media, which in turn ensures that machinery, chemicals, water, gas, and other substances remain within safe operating conditions. Green building engineering is another field whose focus is on regulating indoor thermal conditions (temperature, humidity) under substantial variations of the ambient conditions to provide high-quality living and working environments. In many manufacturing processes, thermal conditioning is also intended to improve the final mechanical (e.g., hardness, toughness, resistance), electrical, or optical properties of the product, the general picture being that high temperature gradients are useful to speed up the process but generally harm the quality of the outcome because of heat transfer inhomogeneities caused by the increased convection by the fluid particles. All such problems fall under the purview of this line of study.

*Corresponding author

Email address: elie.hachem@mines-paristech.fr (E. Hachem)

15 Numerous strategies have been implemented to control fluid mechanical systems (including con-
16 jugate heat transfer systems combining thermal conduction in the solid and convective transfer in
17 the fluid), either open-loop with passive appendices (e.g., end plate, splitter plate, small secondary
18 cylinder, or flexible tail), or open-loop with actuating devices (e.g., plasma actuation, boundary
19 temperatures, steady or unsteady base bleeding, rotation) or closed-loop (e.g. via transverse mo-
20 tion, perturbations of the thermal boundary layer, blowing/suction, rotation, all relying on an
21 appropriate sensing of flow variables). Nonetheless, many of the proposed strategies are trial and
22 error, and therefore require extensive and costly experimental or numerical campaigns. This has
23 motivated the development of analytical methods and numerical algorithms for the optimal control
24 of Navier–Stokes systems [1–3], and the maturing of mathematical methods in flow control and
25 discrete concepts for PDE constrained optimization. Applications to the heat equation [4] and the
26 coupled Navier–Stokes and heat equations [5–8] have also been considered, including fresh devel-
27 opments meant to alter the linear amplification of flow disturbances [9], but the general picture
28 remains that the optimal control of conducting-convecting (possibly radiating) fluids has not been
29 extensively studied.

30 The premise of this research is that the related task of selecting an optimal subset of control
31 parameters can alternatively be assisted by machine learning algorithms. Indeed, the introduc-
32 tion of the back-propagation algorithm [10] has progressively turned Artificial Neural Networks
33 (ANN) into a family of versatile non-parametric tools that can be trained to hierarchically extract
34 informative features from data and to provide qualitative and quantitative modeling predictions.
35 Together with the increased affordability of high-performance computational hardware, this has
36 allowed leveraging the ever-increasing volume of data generated for research and engineering pur-
37 poses into novel insight and actionable information, which in turn has entirely transformed scientific
38 disciplines, such as robotics [11, 12] or image analysis [13]. Owing to the ability of neural networks
39 to handle stiff, large-scale nonlinear problems [14], machine learning algorithms have also been
40 making rapid inroads in fluid mechanics, as a mean to solve the Navier–Stokes equations [15] or
41 to predict closure terms in turbulence models [16]; see also Ref. [17] for an overview of the current
42 developments and opportunities.

43 Neural networks can also be used to solve decision-making problems, which is the purpose of
44 Deep Reinforcement Learning (DRL, where the *deep* terminology generally weighs on the sizable
45 depth of the network), an advanced branch of machine learning. Simply put, a neural network trains
46 in finding out which actions or succession of actions maximize a numerical reward signal, with the
47 possibility for a given action to affect not only the immediate but also the future rewards. Successful
48 applications of DRL range from AlphaGo, the well-known ANN that defeated the top-level human
49 player at the game of Go [18] to the real-world deployment of legged robots [19], to breakthroughs in
50 computer vision (e.g., filtering, or extracting image features) [20] and optimal control problems [21,
51 22]. There is also great potential for applying DRL to fluid mechanics, for which efforts are ongoing
52 but still at an early stage. Sustained commitment from the machine learning community has
53 allowed expanding the scope from computationally inexpensive, low-dimensional reductions of the
54 underlying fluid dynamics [23–25] to complex Navier–Stokes systems [26, 27], with a handful of
55 pioneering studies providing insight into the performance improvements to be delivered in shape
56 optimization [28–30] and flow control [31–37], including recent advances assessing experimentally
57 the effectiveness of reinforcement learning control strategies [38]. The literature on thermal control
58 is even more scarce, as our literature review did not reveal any other study considering DRL-based
59 control of conjugate heat transfer aside from [39], another research effort conducted in the same
60 time frame as the present work that will be discussed further on, plus a few other publications
61 relying on dealing with energy efficiency in civil engineering from low-dimensional thermodynamic
62 models basically unrelated to the equations of fluid dynamics [40, 41].

63 This research assesses the feasibility of using proximal policy optimization (PPO [22]) for control
64 and optimization purposes of conjugate heat transfer systems, as governed by the coupled Navier–
65 Stokes and heat equations. The objective here is to keep shaping the capabilities of the method
66 (PPO is still a relatively newcomer that has quickly emerged as the go-to DRL algorithm due to
67 its data efficiency, simplicity of implementation and reliable performance) and to narrow the gap
68 between DRL and advanced numerical methods for multi-scale, multi-physics computational fluid
69 dynamics (CFD). We investigate more specifically the “degenerate” single-step PPO algorithm
70 introduced in [30] for optimization and open-loop control problems, as the optimal policy to be

71 learnt is then state-independent, and it may be enough for the neural network to get only one
72 attempt per episode at finding the optimal. Several problems of conjugate heat transfer in two and
73 three dimensions are used as testbed to push forward the development of this novel approach, whose
74 potential for reliable black-box optimization of computational fluid dynamics (CFD) systems has
75 been recently assessed for open-loop drag reduction in cylinder flows at Reynolds numbers ranging
76 from a few hundreds to a few ten thousands [42]. To the best of the authors knowledge, this
77 constitutes the first attempt to achieve DRL-based control of conjugate *forced* convection heat
78 transfer, while [39] is the first attempt to achieve DRL control of conjugate *natural* convection
79 heat transfer.

80 The organization is as follows: section 2 outlines the main features of the finite element CFD
81 environment used to compute the numerical reward fed to the neural network, that combines
82 variational multi-scale (VMS) modeling of the governing equations, immerse volume method, and
83 multi-component anisotropic mesh adaptation. The baseline principles and assumptions of DRL
84 and PPO are presented in section 3, together with the specifics of the single-step PPO algorithm.
85 Section 4 revisits the natural convection case of [39] for the purpose of validation and assessment
86 part of the method capabilities. In section 5, DRL is used to control conjugate heat transfer in
87 a model setup of two-dimensional workpiece cooling by impingement of a fluid. An extension to
88 three-dimensional workpieces is proposed in section 6.

89 2. Computational fluid dynamics

90 The focus of this research is on conjugate heat transfer and laminar, incompressible fluid flow
91 problems in two and three-dimensions, for which the conservation of mass, momentum and energy
92 is described by the nonlinear, coupled Navier–Stokes and heat equations

$$\nabla \cdot \mathbf{u} = 0, \quad (1)$$

$$\rho(\partial_t \mathbf{u} + \mathbf{u} \cdot \nabla \mathbf{u}) = \nabla \cdot (-p\mathbf{I} + 2\mu\boldsymbol{\varepsilon}(\mathbf{u})) + \boldsymbol{\psi}, \quad (2)$$

$$\rho c_p(\partial_t T + \mathbf{u} \cdot \nabla T) = \nabla \cdot (\lambda \nabla T) + \chi, \quad (3)$$

93 where \mathbf{u} is the velocity field, p is the pressure, T is the temperature, $\boldsymbol{\varepsilon}(\mathbf{u}) = (\nabla \mathbf{u} + \nabla \mathbf{u}^T)/2$ is the
94 rate of deformation tensor, $\boldsymbol{\psi}$ and χ are source terms (modeling, e.g., buoyancy or radiative heat
95 transfer), and we assume here constant fluid density ρ , dynamic viscosity μ , thermal conductivity
96 λ , and specific heat c_p .

97 2.1. The immersed volume method

98 The numerical modeling of conjugate heat transfer mostly depends upon a heat transfer co-
99 efficient to ensure that the proper amount of heat is exchanged at the fluid/solid interface via
100 thermal boundary conditions. Computing said coefficient is no small task (as it requires solving
101 an inverse problem to assimilate relevant experimental data, which in turn requires such data to
102 be available), and is generally acknowledged to be a limiting issue for practical applications where
103 one must vary, e.g., the shape, number and position of the solid, or the fluid and/or solid material
104 properties. We thus rather use here the immerse volume method (IVM) to combine both the fluid
105 and solid phases into a single fluid with variable material properties. Simply put, we solve equa-
106 tions formally identical to (1)-(3) on a unique computational domain Ω , but with variable density,
107 dynamic viscosity, conductivity, and specific heat, which removes the need for a heat transfer coef-
108 ficient since the amount of heat exchanged at the interface then proceeds solely from the individual
109 material properties on either side of it. In order to ensure numerical accuracy, such an approach
110 must combine three key ingredients, that are briefly reviewed in the next paragraphs: an interface
111 capturing method, anisotropic mesh adaptation to achieve a high-fidelity description of said inter-
112 face, and relevant mixing laws to describe the properties of the composite fluid. One point worth
113 mentioning is that the interface here is static, although the same numerical framework can be used
114 to dynamically track moving interfaces, and thus to encompass the effect of solid displacements.
115 This is because the solid is fixed once an action has been taken by the PPO agent, although not
116 fixed over the course of optimization, as the solid position can very well be the quantity subjected
117 to optimization, as illustrated in section 5.3.4.

118

119 - *Level set method.* The level set approach is used to localize the fluid/solid interface by the zero
 120 iso-value of a smooth function. In practice, a signed distance function ϕ is used to localize the
 121 interface and initialize the material properties on both either side of it, with the convention that
 122 $\phi > 0$ (resp. $\phi < 0$) in the fluid (resp. the solid).
 123

124 - *Anisotropic mesh adaptation.* The interface may intersect arbitrarily the mesh elements if it
 125 is not aligned with the element edges, in which case discontinuous material properties across
 126 the interface can yield oscillations of the numerical solutions. We thus use the anisotropic mesh
 127 adaptation technique presented in [43] to ensure that the material properties are distributed as
 128 accurately and smoothly as possible over the smallest possible thickness around the interface. This
 129 is done computing modified distances from a symmetric positive defined tensor (the metric) whose
 130 eigenvectors define preferential directions along which mesh sizes can be prescribed from the related
 131 eigenvalues. The metric used here is isotropic far from the interface, with mesh size set equal to
 132 h_∞ in all directions, but anisotropic near the interface, with mesh size equal to h_\perp in the direction
 133 normal to the interface, and to h_∞ in the other directions. This is written for an intended thickness
 134 δ as

$$\mathbf{M} = K(\phi)\mathbf{n} \otimes \mathbf{n} + \frac{1}{h_\infty^2}\mathbf{I} \quad \text{with} \quad K(\phi) = \begin{cases} 0 & \text{if } |\phi| \geq \delta/2, \\ \frac{1}{h_\perp^2} - \frac{1}{h_\infty^2} & \text{if } |\phi| < \delta/2, \end{cases} \quad (4)$$

135 where $\mathbf{n} = \nabla\phi/|\nabla\phi|$ is the unit normal to the fluid/solid interface computed from the level set
 136 gradient. A posteriori anisotropic error estimator is then used to minimize the interpolation error
 137 under the constraint of a fixed number of edges in the mesh. A unique metric can be built from
 138 multi-component error vectors [43–46], which is especially relevant for conjugate heat transfer op-
 139 timization, as it allows each learning episode to use an equally accurate mesh adapted from the
 140 velocity vector and magnitude, the temperature field, and the level set.
 141

142 - *Mixing laws.* The composite density, dynamic viscosity and specific heat featured in equations (1)-
 143 (3) are computed as the arithmetic means of the fluid and solid values, for instance the composite
 144 density is

$$\rho = \rho_f H_\epsilon(\phi) + \rho_s(1 - H_\epsilon(\phi)), \quad (5)$$

145 where H_ϵ is the smoothed Heaviside function defined as

$$H_\epsilon(\phi) = \begin{cases} 0 & \text{if } \phi < -\epsilon, \\ \frac{1}{2}\left(1 + \frac{\phi}{\epsilon} + \frac{1}{\pi} \sin\left(\pi \frac{\phi}{\epsilon}\right)\right) & \text{if } |\phi| \leq \epsilon, \\ 1 & \text{if } \phi > \epsilon, \end{cases} \quad (6)$$

146 and ϵ is a regularization parameter proportional to the mesh size in the normal direction to the
 147 interface, set here to $\epsilon = 2h_\perp$. In order to ensure continuity of the heat flux across the interface,
 148 the thermal conductivity is computed as the harmonic mean

$$\frac{1}{\lambda} = \frac{1}{\lambda_f} H_\epsilon(\phi) + \frac{1}{\lambda_s} (1 - H_\epsilon(\phi)), \quad (7)$$

149 as obtained from a steady, no source, one dimensional analysis of the heat flux when the conduc-
 150 tivity varies stepwise from one medium to the next; see [47] for detailed derivation and analysis,
 151 and [48] for proof of the gain in numerical accuracy (with respect to the arithmetic mean model)
 152 by comparison with analytical solutions.

153 2.2. Variational multi-scale approach (VMS)

154 In the context of finite element methods (that remain widely used to simulate engineering
 155 CFD systems due to their ability to handle complex geometries), direct numerical simulation

156 (DNS) solves the weak form of (1)-(3), obtained by integrating by parts the pressure, viscous and
 157 conductive terms, to give

$$(\rho(\partial_t \mathbf{u} + \mathbf{u} \cdot \nabla \mathbf{u}), \mathbf{w}) + (2\mu \boldsymbol{\varepsilon}(\mathbf{u}), \boldsymbol{\varepsilon}(\mathbf{w})) - (p, \nabla \cdot \mathbf{w}) + (\nabla \cdot \mathbf{u}, q) = (\boldsymbol{\psi}, \mathbf{w}), \quad (8)$$

$$(\rho c_p(\partial_t T + \mathbf{u} \cdot \nabla T), s) + (\lambda \nabla T, \nabla s) = (\chi, s), \quad (9)$$

158 where (\cdot, \cdot) is the L^2 inner product on the computational domain, \mathbf{w} , q and s are relevant test
 159 functions for the velocity, pressure and temperature variables, and all fluid properties are those
 160 mixed with the smoothed Heaviside function (6).

161 We use here the variational multi-scale (VMS) approach [49–51] to solve a stabilized formulation
 162 of (8)-(9), which allows circumventing the Babuska–Brezzi condition (that otherwise imposes that
 163 different interpolation orders be used to discretize the velocity and pressure variables, while we
 164 use here simple continuous piecewise linear P_1 elements for all variables) and prevents numerical
 165 instabilities in convection regimes at high Reynolds numbers. We shall not go into the extensive
 166 details about the derivation of the stabilized formulations, for which the reader is referred to [52, 53].
 167 Suffice it to say here that the flow quantities are split into coarse and fine scale components, that
 168 correspond to different levels of resolution. The fine scales are solved in an approximate manner
 169 to allow modeling their effect into the large-scale equations. This gives rise to additional terms in
 170 the right-hand side of (8)-(9), and yields the following weak forms for the large scale

$$\begin{aligned} &(\rho(\partial_t \mathbf{u} + \mathbf{u} \cdot \nabla \mathbf{u}), \mathbf{w}) + (2\mu \boldsymbol{\varepsilon}(\mathbf{u}), \boldsymbol{\varepsilon}(\mathbf{w})) - (p, \nabla \cdot \mathbf{w}) + (\nabla \cdot \mathbf{u}, q) = (\boldsymbol{\psi}, \mathbf{w}) \\ &\quad + \sum_{K \in \mathcal{T}_h} [(\tau_1 \mathcal{R}_M, \mathbf{u} \cdot \nabla \mathbf{w})_K + (\tau_1 \mathcal{R}_M, \nabla q)_K + (\tau_2 \mathcal{R}_C, \nabla \cdot \mathbf{w})_K], \end{aligned} \quad (10)$$

$$\begin{aligned} &(\rho c_p(\partial_t T + \mathbf{u} \cdot \nabla T), s) + (\lambda \nabla T, \nabla s) = (\chi, s) \\ &\quad + \sum_{K \in \mathcal{T}_h} [(\tau_3 \mathcal{R}_T, \mathbf{u} \cdot \nabla s)_K + (\tau_4 \mathcal{R}_T, \zeta \nabla T \cdot \nabla s)_K], \end{aligned} \quad (11)$$

171 where $(\cdot, \cdot)_K$ is the inner product on element K , we denote by $\zeta = \mathbf{u} \cdot \nabla T / \|\nabla T\|^2$ the (normal-
 172 ized) velocity projected along the direction of the temperature gradient, and the \mathcal{R} terms are the
 173 governing equations residuals

$$-\mathcal{R}_C = \nabla \cdot \mathbf{u}, \quad -\mathcal{R}_M = \rho(\partial_t \mathbf{u} + \mathbf{u} \cdot \nabla \mathbf{u}) + \nabla p - \boldsymbol{\psi} \quad -\mathcal{R}_T = \rho c_p(\partial_t T + \mathbf{u} \cdot \nabla T) - \chi, \quad (12)$$

174 whose second derivatives vanish since we use linear interpolation functions. In (10), $\tau_{1,2}$ are ad-hoc
 175 mesh-dependent stabilization parameters defined in [54, 55]. Conversely, in (11), $\tau_{3,4}$ are mesh-
 176 independent stabilization parameters acting both in the direction of the solution and of its gradient,
 177 that proceed from the stabilization of the ubiquitous convection-diffusion-reaction equation [56, 57],
 178 whose definition is given in [58, 59].

179 The governing equations are solved sequentially, i.e., we solve first (10), then use the resulting
 180 fluid velocity to solve (11). All linear systems are preconditioned with a block Jacobi method
 181 supplemented by an incomplete LU factorization, and solved with the GMRES algorithm, with
 182 tolerance threshold set to 10^{-6} for the Navier–Stokes equations, and 10^{-5} for the heat equation.
 183 The time derivatives and convection terms of the Navier–Stokes equations and related VMS source
 184 terms are integrated semi-implicitly using the first-order backward differentiation formula and
 185 Newton–Gregory backward polynomial. The viscous, pressure and divergence terms are treated
 186 implicitly with the backward Euler scheme. Finally, the VMS stabilization terms $\tau_{1,2}$ are treated
 187 explicitly with the forward Euler scheme, which yields

$$\begin{aligned} &\rho \left(\frac{\mathbf{u}^{i+1} - \mathbf{u}^i}{\Delta t} + \mathbf{u}^i \cdot \nabla \mathbf{u}^{i+1} \right), \mathbf{w}) + (2\mu \boldsymbol{\varepsilon}(\mathbf{u}^{i+1}), \boldsymbol{\varepsilon}(\mathbf{w})) - (p^{i+1}, \nabla \cdot \mathbf{w}) + (\nabla \cdot \mathbf{u}^{i+1}, q) = (\boldsymbol{\psi}^i, \mathbf{w}) \\ &\quad + \sum_{K \in \mathcal{T}_h} [(\tau_1^i \mathcal{R}_M^{i+1}, \mathbf{u}^i \cdot \nabla \mathbf{w})_K + (\tau_1^i \mathcal{R}_M^{i+1}, \nabla q)_K + (\tau_2^i \mathcal{R}_C^{i+1}, \nabla \cdot \mathbf{w})_K], \end{aligned} \quad (13)$$

188 with residuals

$$-\mathcal{R}_C^{i+1} = \nabla \cdot \mathbf{u}^{i+1}, \quad -\mathcal{R}_M^{i+1} = \rho \left(\frac{\mathbf{u}^{i+1} - \mathbf{u}^i}{\Delta t} + \mathbf{u}^i \cdot \nabla \mathbf{u}^{i+1} \right) + \nabla p^{i+1} - \boldsymbol{\psi}^i, \quad (14)$$

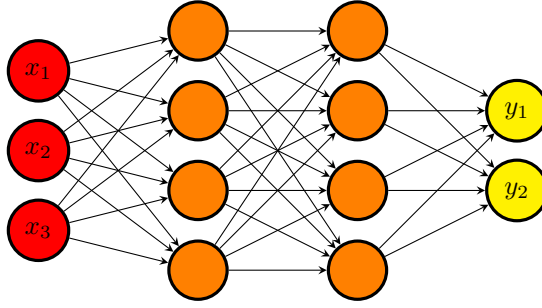


Figure 1: Fully connected neural network with two hidden layers, modeling a mapping from \mathbb{R}^3 to \mathbb{R}^2 .

189 where the superscript i refers to the solution at time $t_i = i\Delta t$. The time derivatives, convection
190 and conduction terms of the heat equation and related VMS source terms are integrated implicitly
191 with the backward Euler scheme (modeling the velocity after \mathbf{u}^{i+1} wherever needed on behalf of
192 the sequential resolution process, although we drop the dependence in the notation to ease the
193 reading). The VMS stabilization terms $\tau_{3,4}$ are treated explicitly with the forward Euler scheme¹,
194 to give

$$\begin{aligned} & (\rho c_p (\frac{T^{i+1} - T^i}{\Delta t} + \mathbf{u}^{i+1} \cdot \nabla T^{i+1}), s) + (\lambda \nabla T^{i+1}, \nabla s) = (\chi^i, s) \\ & + \sum_{K \in \mathcal{T}_h} [(\tau_3^i \mathcal{R}_T^{i+1}, \mathbf{u}^{i+1} \cdot \nabla s)_K + (\tau_4^i \mathcal{R}_T^{i+1}, \zeta^i \nabla T^i \cdot \nabla s)_K], \end{aligned} \quad (15)$$

195 with residual

$$-\mathcal{R}_T^{i+1} = \rho c_p (\frac{T^{i+1} - T^i}{\Delta t} + \mathbf{u}^{i+1} \cdot \nabla T^{i+1}) - \chi^i. \quad (16)$$

196 We solve equations (13)-(15) with an in-house VMS solver whose flexibility, accuracy and reliability
197 is assessed in a series of previous papers to which the reader is referred for further information, see
198 in particular [55, 60] for the detailed mathematical formulation of the IVM in the context of finite
199 element VMS methods. The ability of the IVM to handle the abrupt conductivity change across
200 the fluid/solid interface is documented in [53, 61, 62]. Excellent agreement with reference solutions
201 available from the literature and in-house data obtained enforcing proper thermal conditions at the
202 boundary of body-fitted meshes is reported for several time-dependent conjugate heat transfer test
203 cases (e.g., mixed convection in a plane channel flow, combined convection in square enclosures
204 and conduction/radiation heat transfer, all in two dimensions). Ref. [61] also reports favorable
205 agreement between the IVM and in-house experimental data pertaining to a three-dimensional test
206 case representative of an industrial cooling system, which provides strong evidence of relevance for
207 the intended application.

208 3. Deep reinforcement learning and proximal policy optimization

209 3.1. Neural networks

210 A neural network (NN) is a collection of artificial neurons, i.e., connected computational units
211 that can be trained to arbitrarily well approximate the mapping function between input and output
212 spaces. Each connection provides the output of a neuron as an input to another neuron. Each
213 neuron performs a weighted sum of its inputs, to assign significance to the inputs with regard to the
214 task the algorithm is trying to learn. It then adds a bias to better represent the part of the output
215 that is actually independent of the input. Finally, it feeds an activation function that determines
216 whether and to what extent the computed value should affect the outcome. As sketched in figure 1,

¹That is, with respect to T , but the velocity is modeled after its latest computed approximation \mathbf{u}^{i+1} wherever needed.

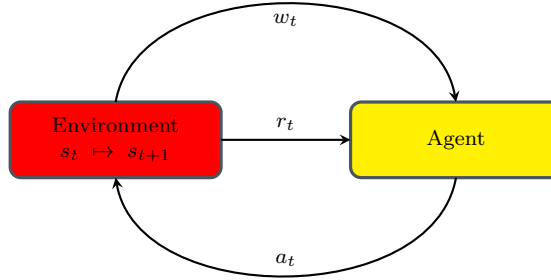


Figure 2: RL agent and its interactions with its environment.

217 a fully connected network is generally organized into layers, with the neurons of one layer being
 218 connected solely to those of the immediately preceding and following layers. The layer that receives
 219 the external data is the input layer, the layer that produces the outcome is the output layer, and
 220 in between them are zero or more hidden layers.

221 The design of an efficient neural network requires a proper optimization of the weights and
 222 biases, together with a relevant nonlinear activation function. The abundant literature available
 223 on this topic points to a relevant network architecture (e.g., type of network, depth, width of each
 224 layer), finely tuned hyper parameters (i.e., parameters whose value cannot be estimated from data,
 225 e.g., optimizer, learning rate, batch size) and a sufficiently large amount of data to learn from as
 226 being the key ingredients for success; see, e.g., Ref. [63] and the references therein.

227 3.2. Deep reinforcement learning

228 Deep reinforcement learning (DRL) is an advanced branch of machine learning in which deep
 229 neural networks train in solving sequential decision-making problems. It is a natural extension of
 230 reinforcement learning (RL), in which an agent (the neural network) is taught how to behave in an
 231 environment by taking actions and by receiving feedback from it under the form of a reward (to
 232 measure how good or bad the action was) and information (to gauge how the action has affected the
 233 environment). This can be formulated as a Markov Decision Process, for which a typical execution
 234 goes as follows (see also figure 2):

- 235 • assume the environment is in state $s_t \in \mathcal{S}$ at iteration t , where \mathcal{S} is a set of states,
- 236 • the agent uses w_t , an observation of the current environment state (and possibly a partial
 237 subset of s_t) to take action $a_t \in \mathcal{A}$, where \mathcal{A} is a set of actions,
- 238 • the environment reacts to the action and transitions from s_t to state $s_{t+1} \in \mathcal{S}$,
- 239 • the agent is fed with a reward $r_t \in \mathcal{R}$, where \mathcal{R} is a set of rewards, and a new observation
 240 w_{t+1} ,

241 This repeats until some termination state is reached, the succession of states and actions defining
 242 a trajectory $\tau = (s_0, a_0, s_1, a_1, \dots)$. In any given state, the objective of the agent is to determine
 243 the action maximizing its cumulative reward over an episode, i.e., over one instance of the scenario
 244 in which the agent takes actions. Most often, the quantity of interest is the discounted cumulative
 245 reward along a trajectory defined as

$$R(\tau) = \sum_{t=0}^T \gamma^t r_t, \quad (17)$$

246 where T is the horizon of the trajectory, and $\gamma \in [0, 1]$ is a discount factor that weighs the relative
 247 importance of present and future rewards (the agent being short-sighted in the limit where $\gamma \rightarrow 0$,
 248 since it then cares solely about the first reward, and far-sighted in the limit where $\gamma \rightarrow 1$, since it
 249 then cares equally about all rewards).

250 There exist two main types of RL algorithms, namely model-based methods, in which the
 251 agent tries to build a model of how the environment works to make predictions about what the
 252 next state and reward will be before taking any action, and model-free methods, in which the agent
 253 conversely interacts with the environment without trying to understand it, and are prominent in
 254 the DRL community. Another important distinction to be made within model-free algorithms
 255 is that between value-based methods, in which the agent learns to predict the future reward of
 256 taking an action when provided a given state, then selects the maximum action based on these
 257 estimates, and policy-based methods, in which it optimizes the expected reward of a decision policy
 258 mapping states to actions. Many of the most successful algorithms in DRL (including proximal
 259 policy optimization, whose assessment for flow control and optimization purposes is the primary
 260 motivation for this research) proceed from policy gradient methods, in which gradient ascent is
 261 used to optimize a parameterized policy with respect to the expected return, as further explained
 262 in the next section. The reader interested in a more thorough introduction to the zoology of RL
 263 methods (together with their respective pros and cons) is referred to Ref. [64].

264 3.3. From policy methods to Proximal policy optimization

265 This section intended for the non-specialist reader briefly reviews the basic principles and as-
 266 sumptions of policy gradient methods, together with the various steps taken for improvement.
 267

268 - *Policy methods.* A policy method maximizes the expected discounted cumulative reward of a
 269 decision policy mapping states to actions. It resorts not to a value function, but to a probability
 270 distribution over actions given states, that fully defines the behavior of the agent. Since policies
 271 are most often stochastic, the following notations are introduced:

- 272 • $\pi(s, a)$ is the probability of taking action a in state s under policy π ,
- 273 • $Q^\pi(s, a)$ is the expected value of the return of the policy after taking action a in state s (also
 274 termed state-action value function or Q-function)

$$Q^\pi(s, a) = \mathbb{E}_\pi [R(\tau)|s, a], \quad (18)$$

275 where we use \mathbb{E}_π for the expected value \mathbb{E} under policy π .

- 276 • $V^\pi(s)$ is the expected value of the return of the policy in state s (also termed value function
 277 or V-function)

$$V^\pi(s) = \mathbb{E}_\pi [R(\tau)|s]. \quad (19)$$

278 The V and Q functions are therefore such that

$$V^\pi(s) = \sum_a \pi(s, a) Q^\pi(s, a), \quad (20)$$

279 so $V^\pi(s)$ can also be understood as the probability-weighted average of discounted cumulated
 280 rewards over all possible actions in state s .

281 - *Policy gradient methods.* A policy gradient method aims at optimizing a parametrized policy
 282 π_θ , where θ denotes the free parameters whose value can be learnt from data (as opposed to the
 283 hyper parameters). In practice, one defines an objective function based on the expected discounted
 284 cumulative reward

$$J(\theta) = \mathbb{E}_{\pi_\theta} [R(\tau)], \quad (21)$$

285 and seeks the parameterization θ^* maximizing $J(\theta)$, hence such that

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\pi_\theta} [R(\tau)], \quad (22)$$

286 which can be done on paper by plugging an estimator of the policy gradient $\nabla_\theta J(\theta)$ into a gradient
 287 ascent algorithm. This is no small task as one is looking for the gradient with respect to the policy

288 parameters, in a context where the effects of policy changes on the state distribution are unknown
 289 (since modifying the policy will most likely modify the set of visited states, which will in turn affect
 290 performance in some indefinite manner). One commonly used estimator, derived in [64] using the
 291 log-probability trick, reads

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log(\pi_{\theta}(s_t, a_t)) R(\tau) \right] \sim \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log(\pi_{\theta}(s_t, a_t)) \widehat{A}^{\pi}(s_t, a_t) \right], \quad (23)$$

292 where \widehat{A}^{π} is some biased estimator (here its normalization to zero mean and unit variance) of the
 293 advantage function

$$A^{\pi}(s, a) = Q^{\pi}(s, a) - V^{\pi}(s), \quad (24)$$

294 that measures the improvement (if $A^{\pi} > 0$, otherwise the lack thereof) associated with taking
 295 action a in state s compared to taking the average over all possible actions. This is because the
 296 value function does not depend on θ , so taking it off changes neither the expected value, nor the
 297 gradient, but it does reduce the variance, and speeds up the training. Furthermore, when the
 298 policy π_{θ} is represented by a neural network (in which case θ simply denotes the network weights
 299 and biases to be optimized), the focus is rather on the policy loss defined as

$$L(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^T \log(\pi_{\theta}(a_t|s_t)) \widehat{A}(s_t, a_t) \right], \quad (25)$$

300 whose gradient is equal to the (approximated) policy gradient (23) (since the gradient operator acts
 301 only on the log-policy term, not on the advantage) and is computed with respect to each weight
 302 and bias by the chain rule, one layer at the time, using the back-propagation algorithm [10].
 303

304 - *Trust regions.* The performance of policy gradient methods is hurt by the high sensitivity to the
 305 learning rate, i.e., the size of the step to be taken in the gradient direction. Indeed, small learning
 306 rates are detrimental to learning, but large learning rates can lead to a performance collapse if the
 307 agent falls off the cliff and restarts from a poorly performing state with a locally bad policy. This
 308 is all the more harmful as the learning rate cannot be tuned locally, meaning that an above average
 309 learning rate will speed up learning in some regions of the parameter space where the policy loss
 310 is relatively flat, but will possibly trigger an exploding policy update in other regions exhibiting
 311 sharper variations. One way to ensure continuous improvement is by imposing a trust region con-
 312 straint to limit the difference between the current and updated policies, which can be done by
 313 determining first a maximum step size relevant for exploration, then by locating the optimal point
 314 within this trust region. We will not dwell on the intricate details of the many algorithms developed
 315 to solve such trust region optimization problems, e.g., natural policy gradient (NPG [65]), or trust
 316 region policy optimization (TRPO [66]). Suffice it to say that they use the minorize-maximization
 317 algorithm to maximize iteratively a surrogate policy loss (i.e. a lower bound approximating locally
 318 the actual loss at the current policy), but are difficult to implement and can be computationally
 319 expensive, as they rely on an estimate of the second-order gradient of the policy log probability.
 320

321 - *Proximal policy optimization.* Proximal policy optimization (PPO) is another approach with
 322 simple and effective heuristics, that uses a probability ratio between the two policies to maximize
 323 improvement without the risk of performance collapse [22]. The focus here is on the PPO-clip
 324 algorithm², that optimizes the surrogate loss

$$L(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\min \left(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)}, g(\epsilon, \widehat{A}^{\pi}(s, a)) \right) \widehat{A}^{\pi}(s, a) \right], \quad (26)$$

²There is also a PPO-Penalty variant which uses a penalization on the average Kullback–Leibler divergence between the current and new policies, but PPO-clip performs better in practice.

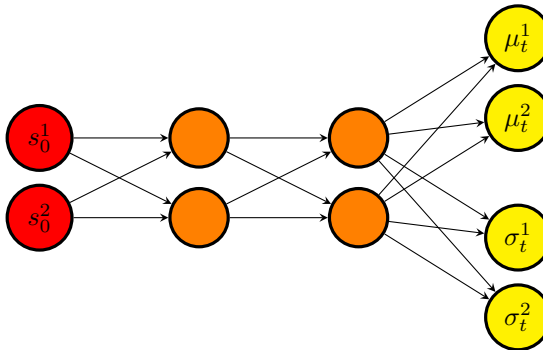


Figure 3: Agent network example used to map states to policy. The input state \mathbf{s}_0 , here of size 2, is mapped to a mean $\boldsymbol{\mu}$ and a standard deviation $\boldsymbol{\sigma}$ vectors, each of size 2. All activation functions are ReLU, except for that of the last layer, which are linear for the $\boldsymbol{\mu}$ output, and softplus for the $\boldsymbol{\sigma}$ output. Orthogonal weights initialization is used throughout the network.

325 where

$$g(\epsilon, A) = \begin{cases} 1 + \epsilon & A \geq 0, \\ 1 - \epsilon & A < 0, \end{cases} \quad (27)$$

326 and $\epsilon \in [0.1, 0.3]$ is the clipping range, a small hyper parameter defining how far away the new policy
 327 is allowed to go from the old. The general picture is that a positive (resp. negative) advantage
 328 increases (resp. decreases) the probability of taking action a in state s , but always by a proportion
 329 smaller than ϵ , otherwise the min kicks in (26) and its argument hits a ceiling of $1 + \epsilon$ (resp. a
 330 floor of $1 - \epsilon$). This prevents stepping too far away from the current policy, and ensures that the
 331 new policy will behave similarly.

332 There exist more sophisticated PPO algorithms (e.g., Trust region PPO [67], that determines
 333 first a maximum step size relevant for exploration, then adaptively adjusts the clipping range to
 334 find the optimal within this trust region), but standard PPO has simple and effective heuristics.
 335 Namely, it is computationally inexpensive, easy to implement (as only the first-order gradient of
 336 the policy log probability is needed to calculate the clipped surrogate), and remains regarded as
 337 one of the most successful RL algorithms, achieving state-of-the-art performance across a wide
 338 range of challenging tasks.

339 3.4. Single-step PPO

340 We now come to single-step PPO, a “degenerate” version of PPO introduced in [30] and intended
 341 for situations where the optimal policy to be learnt by the neural network is state-independent, as
 342 is notably the case in optimization and open-loop control problems (closed-loop control problems
 343 conversely require state-dependent policies for which standard PPO is best suited). The main
 344 difference between standard and single-step PPO can be summed up as follows: where standard
 345 PPO seeks the optimal set of actions a^* yielding the largest possible reward, single-step PPO seeks
 346 the optimal mapping f_{θ^*} such that $a^* = f_{\theta^*}(s_0)$, where θ denotes the network free parameters
 347 and s_0 is some input state (usually a vector of zeros) consistently fed to the agent for the optimal
 348 policy to eventually embody the transformation from s_0 to a^* . The agent initially implements a
 349 random state-action mapping f_{θ_0} from s_0 to an initial policy determined by the free parameters
 350 initialization θ_0 , after which it gets only one attempt per learning episode at finding the optimal
 351 (i.e., it interacts with the environment only once per episode). This is illustrated in figure 4 showing
 352 the agent draw a population of actions $a_t = f_{\theta_t}(s_0)$ from the current policy, and being returned
 353 incentives from the associated rewards to update the free parameters for the next population of
 354 actions $a_{t+1} = f_{\theta_{t+1}}(s_0)$ to yield larger rewards.

355 In practice, the agent outputs a policy parameterized by the mean and variance of the proba-
 356 bility density function of a d -dimensional multivariate normal distribution, with d the dimension
 357 of the action required by the environment. Actions drawn in $[-1, 1]^d$ are then mapped into rele-
 358 vant physical ranges, a step deferred to the environment as being problem-specific. The resolution

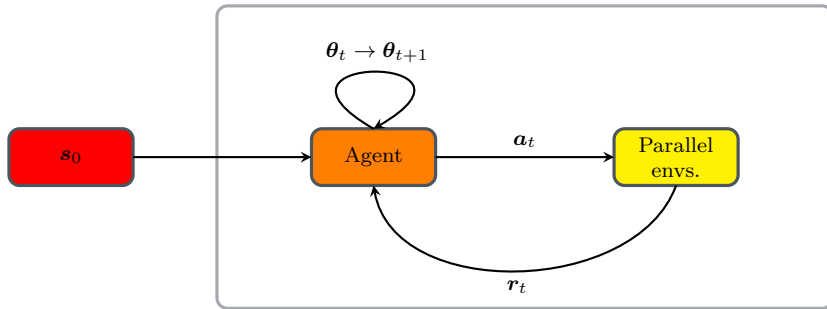


Figure 4: Action loop for single-step PPO. At each episode, the input state s_0 is provided to the agent, which in turn provides n actions to n parallel environments. The latter return n rewards, that evaluate the quality of each action taken. Once all the rewards are collected, an update of the agent parameters is made using the PPO loss (26).

359 essentially follows the process described in section 3.3, only a normalized averaged reward substi-
 360 tutes for the advantage function. This is because classical PPO is actor-critic, i.e., it improves the
 361 learning performance by updating two different networks, a first one called actor that controls the
 362 actions taken by the agent, and a second one called critic, that learns to estimate the advantage
 363 from the value function as

$$A(s_t, a_t) = r_t + \gamma V(s_{t+1}) - V(s_t). \quad (28)$$

364 In single-step PPO, the trajectory consists of a single state-action pair, so the discount factor can
 365 be set to $\gamma = 1$ with no loss of generality. In return, the advantage reduces to the whitened reward
 366 since the two rightmost terms cancel each other out in (28). This means that the approach can do
 367 without the value-function evaluations of the critic network, i.e., it is not actually actor-critic.

368 3.5. Numerical implementation

369 The present workflow relies on the online PPO implementation of Stable Baselines, a toolset
 370 of reinforcement learning algorithms dedicated to the research community and industry [68], for
 371 which a custom OpenAI environment has been designed using the Gym library [69]. Hyperbolic
 372 tangent is used as default activation function. The instant reward r_t used to train the neural
 373 network is simply the quantity subjected to optimization (modulo a plus or minus sign to tackle
 374 both maximization and minimization problems). A moving average reward is also computed on the
 375 fly as the sliding average over the 100 latest values of r_t (or the whole sample if it has insufficient
 376 size). All other relevant hyper parameters are documented in the next sections, with the exception
 377 of the discount factor (set to $\gamma = 1$).

378 In practice, actions are distributed to multiple environments running in parallel, each of which
 379 executes a self-contained MPI-parallel CFD simulation and feeds data to the DRL algorithm (hence,
 380 two levels of parallelism related to the environment and the computing architecture). The algorithm
 381 waits for the simulations running in all parallel environments to be completed, then shuffles and
 382 splits the rewards data set collected from all environments into several buffers (or mini-batches)
 383 used sequentially to compute the loss and perform a network update. The process repeats for
 384 several epochs, i.e., several full passes of the training algorithm over the entire data set (so the
 385 policy network ends up being trained on samples generated by older policies, which is customary in
 386 standard PPO operation). This simple parallelization technique is key to use DRL in the context
 387 of CFD applications, as a sufficient number of actions drawn from the current policy must be
 388 evaluated to accurately estimate the policy gradient. This comes at the expense of computing
 389 the same amount of reward evaluations, and yields a substantial computational cost for high-
 390 dimensional fluid dynamics problems (typically from a few tens to several thousand hours for the
 391 steady-state optimization problems considered herein). In the same vein, it should be noted that
 392 the common practice in DRL studies to gain insight into the performances of the selected algorithm
 393 by averaging results over multiple independent training runs with different random seeds is not
 394 tractable, as it would trigger a prohibitively large computational burden. The same random seeds
 395 have thus been deliberately used over the whole course of study to ensure a minimal level of
 396 performance comparison between cases.

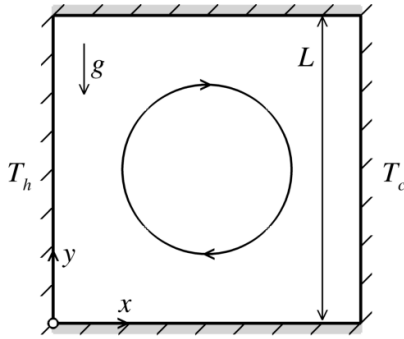


Figure 5: Schematic of the two-dimensional Rayleigh-Bénard set-up.

4. Control of natural convection in 2-D closed cavity

4.1. Case description

We address first the control of natural convection in the two-dimensional differentially heated square cavity schematically illustrated in figure 5(a). This is a widely studied benchmark system for thermally-driven flows, relevant in nature and technical applications (e.g., ocean and atmospheric convection, materials processing, metallurgy), that is thus suitable to validate and compare numerical solution algorithms while enriching the knowledge base for future projects in this field. A Cartesian coordinate system is used with origin at the lower-left edge, horizontal x -axis, and vertical y -axis. The cavity has side L , its top and bottom horizontal walls are perfectly insulated from the outside, and the vertical sidewalls are isothermal. Namely, the right sidewall is kept at a constant, homogeneous “cold” temperature T_c , and the left sidewall is entirely controllable via a constant in time, varying in space “hot” distribution $T_h(y)$ such that

$$\langle T_h \rangle > T_c, \quad (29)$$

where the brackets denote the average over space (here over the vertical position along the sidewall).

In the following, we neglect radiative heat transfer ($\chi = 0$) and consider a Boussinesq system driven by buoyancy, hence

$$\boldsymbol{\psi} = \rho_0 \beta (T - T_c) g \mathbf{e}_y, \quad (30)$$

where \mathbf{g} is the gravitational acceleration parallel to the sidewalls, β is the thermal expansion coefficient, and we use the cold sidewall temperature as Boussinesq reference temperature. By doing so, the pressure featured in the momentum equation (2) and related weak forms must be understood as the pressure correction representing the deviation from hydrostatic equilibrium. The governing equations are solved with no-slip conditions $\mathbf{u} = \mathbf{0}$ on $\partial\Omega$ and temperature boundary conditions

$$\partial_y T(x, 0, t) = \partial_y T(x, L, t) = 0, \quad T(0, y, t) = \langle T_h \rangle + \tilde{T}_h(y), \quad T(L, y, t) = T_c, \quad (31)$$

where \tilde{T}_h is a zero-mean (in the sense of the average over space) distribution of hot temperature fluctuations subjected to optimization, whose magnitude is bounded by some constant ΔT_{max} according to

$$|\tilde{T}_h(y)| \leq \Delta T_{max}, \quad (32)$$

to avoid extreme and nonphysical temperature gradients. All results are made non-dimensional using the cavity side, the heat conductivity time, and the well-defined, constant in time difference between the averaged sidewall temperatures. The retained fluid properties yield values of the Rayleigh and Prandtl numbers

$$\text{Ra} = \frac{g\beta(\langle T_h \rangle - T_c)L^3}{\nu\alpha} = 10^4, \quad \text{Pr} = \frac{\nu}{\alpha} = 0.71, \quad (33)$$

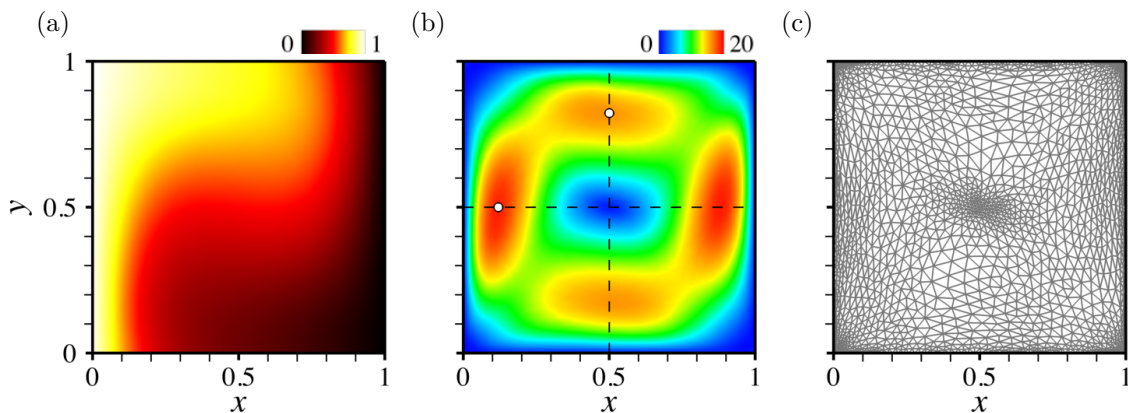


Figure 6: Iso-contours of the uncontrolled steady state (a) temperature and (b) velocity magnitude. (c) Adapted mesh. The circle symbols in (b) mark the positions of the maximum horizontal and vertical velocity along the centerlines reported in table 1.

	Present	Ref. [70]	Ref. [71]	Ref. [72]	Ref. [73]	Ref. [74]
Nu	2.267	2.238	2.245	2.201	2.245	2.245
$\max u(0.5, y)$	16.048	16.178	16.179	–	16.262	16.178
y_{max}	0.823	0.823	0.824	0.832	0.818	0.827
$\max v(x, 0.5)$	19.067	19.617	19.619	–	19.717	19.633
x_{max}	0.120	0.119	0.121	0.113	0.119	0.123

Table 1: Comparison of the present numerical results in the absence of control with reference benchmark solutions from the literature.

425 where $\alpha = \lambda/(\rho c_p)$ is the thermal diffusivity.

426 In order to assess the accuracy of the numerical framework, the uncontrolled solution has
427 been computed by performing 60 iterations with time step $\Delta t = 0.5$ to march the initial solution
428 (consisting of zero velocity and uniform temperature, except at the hot sidewall) to steady state.
429 At each time step, an initially isotropic mesh is adapted under the constraint of a fixed number
430 of elements $n_{el} = 4000$ using a multiple-component criterion featuring velocity and temperature,
431 but no level-set. This is because the case is heat transfer but not conjugate heat transfer, as the
432 solid is solely at the boundary $\partial\Omega$ of the computational domain, where either the temperature is
433 known, or the heat flux is zero. It is thus implemented without the IVM and without a level set
434 (although accurate IVM numerical solutions have been obtained in [53] using thick sidewalls with
435 high thermal conductivity). The solution shown in figure 6(a,b) features a centered roll confined
436 by the cavity walls, consistently with the fact that Ra exceeds the critical value $\text{Ra}_c \sim 920$ for the
437 onset of convection (as extrapolated from the near-critical benchmark data in [70]) by one order
438 of magnitude, and heat transfer is thus driven by both conduction and convection. This shows in
439 the Nusselt number, i.e., the non-dimensional temperature gradient averaged over the hot sidewall

$$\text{Nu} = -\langle \partial_x T \rangle, \quad (34)$$

440 whose present value $\text{Nu} = 2.27$ (as computed from 68 points uniformly distributed along the
441 sidewall) exceeds that $\text{Nu} = 1$ of the purely conductive solution, and exhibits excellent agreement
442 with benchmark results from the literature. This is evidenced in table 1 where we also report the
443 magnitude and position of the maximum horizontal velocity u (resp. the vertical velocity v) along
444 the vertical centerline (resp. the horizontal centerline). The corresponding adapted mesh shown in
445 figure 6(c) stresses that all boundary layers are sharply captured via extremely stretched elements,
446 and that the adaptation strategy yields refined meshes near high temperature gradients and close
447 to the side walls. Note however, the mesh refinement is not only along the boundary layers but also
448 close to the recirculation regions near the cavity center, while the elements in-between are coarse
449 and essentially isotropic.

4.2. Control

The question now being raised is whether DRL can be used to find a distribution of temperature fluctuations \tilde{T}_h capable of alleviating convective heat transfer. To do so, we follow [39] and train a DRL agent in selecting piece-wise constant temperature distributions over n_s identical segments, each of which allows only two pre-determined states referred to as hot or cold. This is intended to reduce the complexity and the computational resources, as large/continuous action spaces are known to be challenging for the convergence of RL methods [28, 75]. Simply put, the network action output consists of n_s values $\hat{T}_{hk \in \{1 \dots n_s\}} = \pm \Delta T_{max}$, mapped into the actual fluctuations according to

$$\tilde{T}_{hk} = \frac{\hat{T}_{hk} - \langle \hat{T}_{hk} \rangle}{\max_l \left\{ 1, \frac{|\hat{T}_{hl} - \langle \hat{T}_{hl} \rangle|}{\Delta T_{max}} \right\}}, \quad (35)$$

to fulfill the zero-mean and upper bound constraints.³ Ultimately, the agent receives the reward $r_t = -\text{Nu}$ to minimize the space averaged heat flux at the hot sidewall.

All results reported herein are for $\Delta T_{max} = 0.75$ (so the hot temperature varies in the range $[0.25; 1.75]$) and $n_s = 10$ segments, as [39] report that $n_s = 20$ was computationally too demanding for their case, and that $n_s = 5$ yielded poor control efficiency. The agent is a fully-connected network with two hidden layers, each holding 2 neurons. The resolution process uses 8 environments and 2 steps mini-batches to update the network for 32 epochs, with learning rate 5×10^{-3} , and PPO loss clipping range $\epsilon = 0.2$.

4.3. Results

For this case, 120 episodes have been run, each of which follows the exact same procedure as above and performs 60 iterations with time step $\Delta t = 0.5$ to march the zero-initial condition to steady state. This represents 960 simulations, each of which is performed on 4 cores and lasts 20s, hence 5h of total CPU cost. We present in figure 7 representative iso-contours of the steady-state temperature and velocity magnitude computed over the course of the optimization. The latter exhibit strong temperature gradients at the hot sidewall, together with a robust steady roll-shaped pattern accompanied by a small corner eddy at the upper-left edge of the cavity, whose size and position depends on the specifics of the temperature distribution. The corresponding meshes are displayed in figure 7(c) to stress the ability of the adaptation procedure to handle well the anisotropy of the solution caused by the intrinsic flow dynamics and the discontinuous boundary conditions.

We show in figure 8 the evolution of the controlled averaged Nusselt number, whose moving average decreases monotonically and reaches a plateau after about 90 episodes, although we notice that sub-optimal distributions keep being explored occasionally. The optimal computed by averaging over the 10 latest episodes (hence the 800 latest instant values) is $\langle \text{Nu} \rangle^* \sim 0.57$, with variations ± 0.01 computed from the root-mean-square of the moving average over the same interval (which is a simple yet robust criterion to assess qualitatively convergence a posteriori). Interestingly, the optimized Nusselt number is almost twice as small as the purely conductive value ($\text{Nu} = 1$), meaning that the approach successfully alleviates the heat transfer enhancement generated by the onset of convection, although it does not alleviate convection itself, as evidenced by the consistent roll-shaped pattern in figure 9. Similar results are reported in [39], for a different set-up in which the horizontal cavity walls are isothermal and control is applied at the bottom at the cavity (hence a different physics because of buoyancy), albeit with lower numerical and control efficiency since the authors report an optimal Nusselt number $\text{Nu} \sim 1$ using up to 512 DRL environments with learning rate of 2.5×10^{-4} . The reason for such discrepancies probably lies in different ways of achieving and assessing control, as we use single-step PPO to optimize the steady-state Nusselt

³Another possible approach would have been to penalize the reward passed to the DRL for those temperature distributions deemed non-admissible (either because the average temperature is non-zero or the temperature magnitude is beyond the threshold). However, this would have made returning admissible solutions part of the tasks the network is trained on (not to mention that non-zero average temperatures amount to a change in the Rayleigh number), which would likely have slowed down learning substantially.

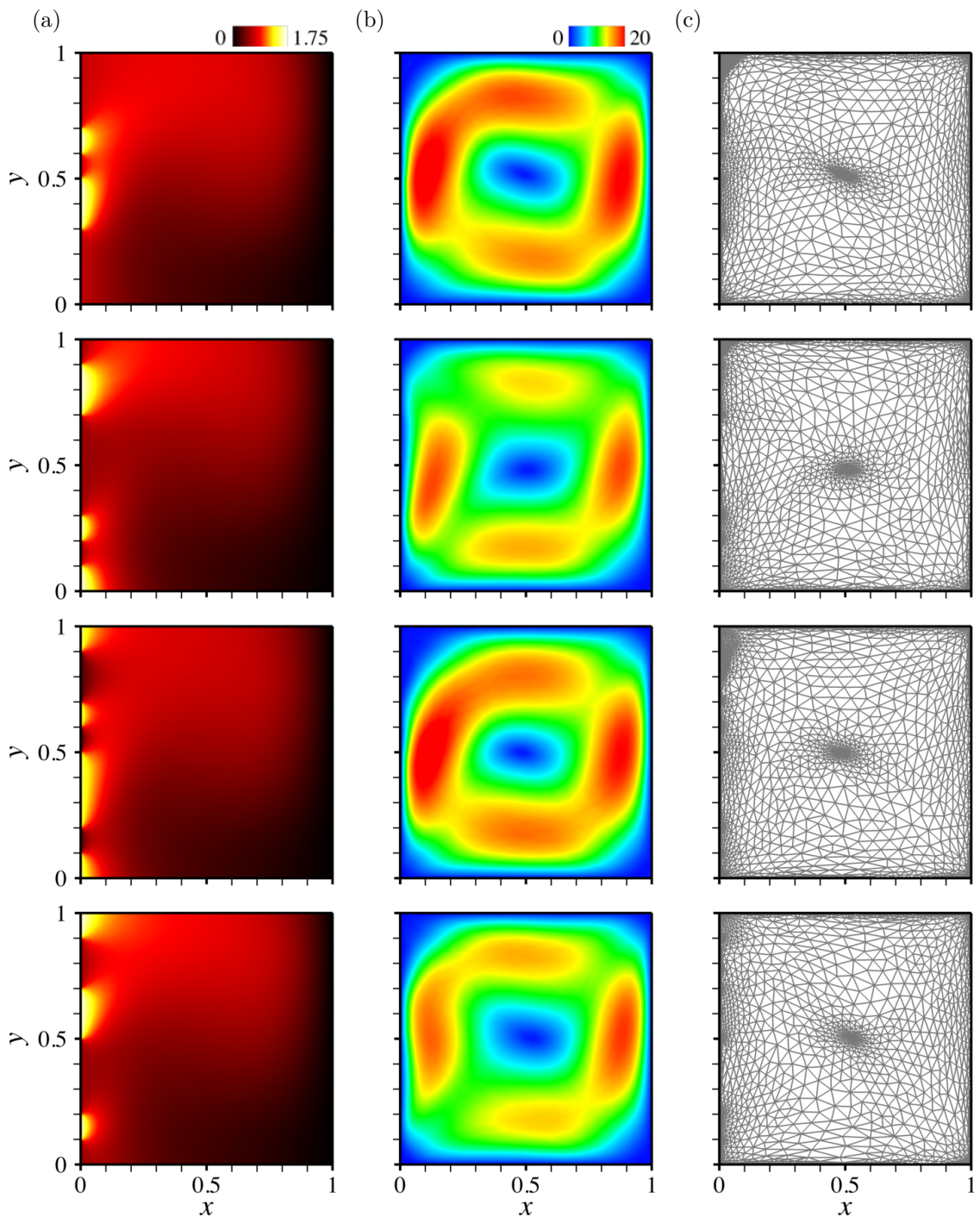


Figure 7: (a,b) Steady-state (a) temperature and (b) velocity magnitude against zero-mean temperature distributed at the left sidewall. (c) Adapted meshes.

494 number via a time-independent control, which requires choosing a sidewall temperature, marching
 495 the controlled solution to steady state, then computing the reward. The problem considered in [39]
 496 is more intricate, as classical PPO is used to optimize the reward accumulated over time via a
 497 time-dependent control temperature updated with a certain period scaling with the convection
 498 time in the cavity (the so-determined optimal control being ultimately time-independent for the
 499 considered value of Ra, but truly time-dependent for Rayleigh numbers above $\sim 10^5$).

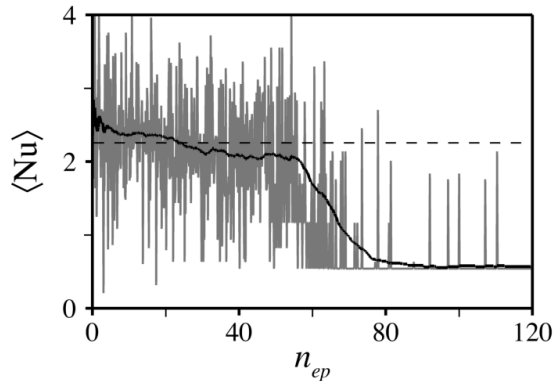


Figure 8: Evolution per learning episode of the instant (in grey) and moving average (in black) Nusselt number. The horizontal dashed line marks the uncontrolled value.

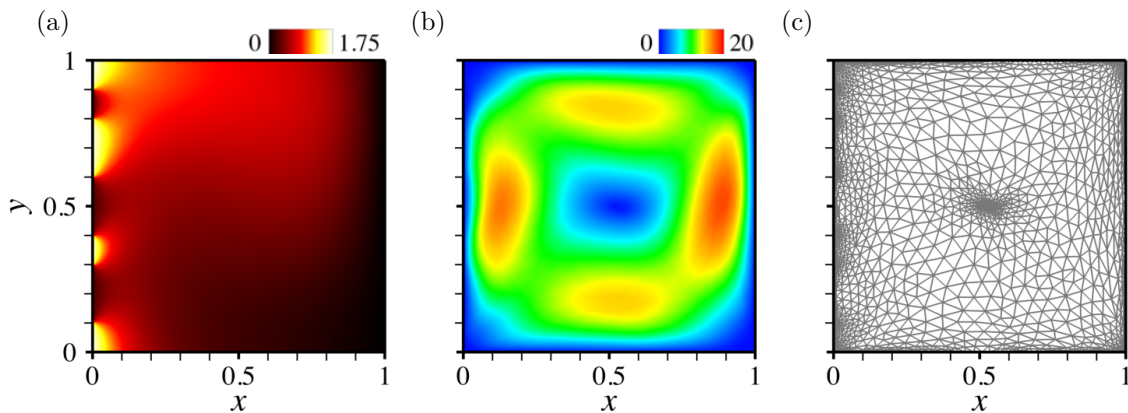


Figure 9: (a,b) Steady-state (a) temperature and (b) velocity magnitude for the optimal zero-mean temperature distribution. (c) Adapted mesh.

500 5. Control of forced convection in 2-D open cavity

501 5.1. Case description

502 This second test case addresses the control of actual conjugate heat transfer in a model setup
 503 for the cooling of a hot solid by impingement of a fluid; see figure 10(a). A Cartesian coordinate
 504 system is used with origin at the center of mass of the solid, horizontal x -axis, and vertical y -axis.
 505 The solid has rectangular shape with height h and aspect ratio 2:1, and is initially at the hot
 506 temperature T_h . It is fixed at the center of a rectangular cavity with height H and aspect ratio
 507 4:1, whose walls are isothermal and kept at temperature T_w . The top cavity side is flush with n_j
 508 identical holes of width e_i whose distribution is subjected to optimization, each of which models
 509 the exit plane of an injector blowing cold air at velocity V_i and temperature T_c , and is identified
 510 by the horizontal position of its center $x_{k \in \{1 \dots n_j\}}$. Hot air is released through the cavity sidewalls,
 511 blown with two identical exhaust areas of height e_o , and identified by the vertical position of their
 512 center $(e_0 - H)/2$.

513 For this case, both buoyancy and radiative heat transfer are neglected ($\psi = \mathbf{0}$ and $\chi = 0$),
 514 meaning that temperature evolves as a passive scalar, similar to the mass fraction of a reactant
 515 in a chemical reaction. All relevant parameters are provided in Table 2, including the material
 516 properties used to model the composite fluid, that yield fluid values of the Reynolds and Prandtl
 517 numbers

$$\text{Re} = \frac{\rho V_i e}{\mu} = 200, \quad \text{Pr} = 2. \quad (36)$$

518 Note the very high value of the solid to fluid viscosity ratio, that ensures that the velocity is zero
 519 in the solid domain and that the no-slip interface condition is satisfied. By doing so, the convective

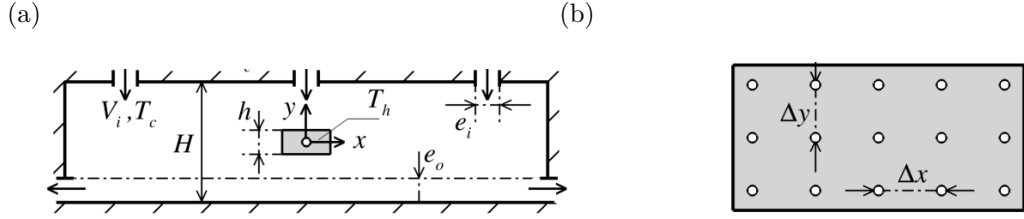


Figure 10: (a) Schematic of the 2-D forced convection set-up. (b) Sensors positions in the solid domain.

H	h	e_i	e_0	V_i	T_w	T_c	T_h	μ	ρ	λ	c_p	
1	0.2	0.2	0.2	1	10	10	150	0.001	1	0.5	1000	Fluid
								1000	100	15	300	Solid

Table 2: Numerical parameters used in the 2-D forced convection problem. All values in SI units, with the exception of temperatures given in Celsius.

520 terms drop out in the energy equation, that reduces to the pure conduction equation for the solid.
 521 The governing equations are solved with no-slip isothermal conditions $\mathbf{u} = \mathbf{0}$ and $T = T_w$ on $\partial\Omega$,
 522 except at the injection exit planes ($\mathbf{u} = -V_i \mathbf{e}_y$, $T = T_c$), and at the exhaust areas, where a zero-
 523 pressure outflow condition is imposed ($p = \partial_x u = \partial_x T = 0$). No thermal condition is imposed at
 524 the interface, where heat exchange is implicitly driven by the difference in the individual material
 525 properties.

5.2. Control

526 The quantity being optimized is the distribution of the injectors center positions $x_{k \in \{1 \dots n_j\}}$.
 527 Several control strategies are assessed in the following, whose ability to manage increasing design
 528 complexity translates into less constrained operation when it comes to optimizing a practically
 529 meaningful device. In practice, each injector is forced to sit in an interval $[x_k^-; x_k^+]$ whose edge
 530 values are determined beforehand or recomputed on the fly (depending on the control strategy),
 531 and bounded according to
 532

$$|x_k^\pm| \leq x_m, \quad (37)$$

533 where we set $x_m = 2H - 0.75e_i$ to avoid numerical issues at the upper cavity edges. The network
 534 action output therefore consists of n_j values $\hat{x} \in [-1; 1]$, mapped into the actual positions according
 535 to

$$x_k = \frac{x_k^+(\hat{x}_k + 1) - x_k^-(\hat{x}_k - 1)}{2}. \quad (38)$$

536 In order to compute the reward passed to the DRL, we distribute uniformly 15 probes in the
 537 solid domain, into $n_x = 5$ columns and $n_y = 3$ rows with resolutions $\Delta x = 0.09$ and $\Delta y = 0.075$,
 538 respectively; see figure 10(b). Selected tests have been carried out to check that the outcome
 539 of the learning process does not change using $n_y = 5$ rows of $n_x = 5$ probes (not shown here).
 540 The magnitude of the tangential heat flux is estimated by averaging the norm of the temperature
 541 gradient over all columns and rows, i.e., i -th column (resp. the j -th row) as

$$\langle \|\nabla_{\parallel} T\| \rangle_i = \frac{2}{n_y - 1} \left| \sum_{j \neq 0} \text{sgn}(j) \|\nabla T\|_{ij} \right|, \quad \langle \|\nabla_{\parallel} T\| \rangle_j = \frac{2}{n_x - 1} \left| \sum_{i \neq 0} \text{sgn}(i) \|\nabla T\|_{ij} \right|, \quad (39)$$

542 where subscripts i , j and ij denote quantities evaluated at $x = i\Delta x$, $y = j\Delta y$ and $(x, y) =$
 543 $(i\Delta x, j\Delta y)$, respectively, and symmetrical numbering is used for the center probe to sit at the
 544 intersection of the zero-th column and row. The numerical reward $r_t = -\langle \|\nabla_{\parallel} T\| \rangle$ fed to the DRL
 545 agent deduces ultimately by averaging over all rows and columns, to give

$$\langle \|\nabla_{\parallel} T\| \rangle = \frac{1}{n_x + n_y} \sum_{i,j} \langle \|\nabla_{\parallel} T\| \rangle_i + \langle \|\nabla_{\parallel} T\| \rangle_j, \quad (40)$$

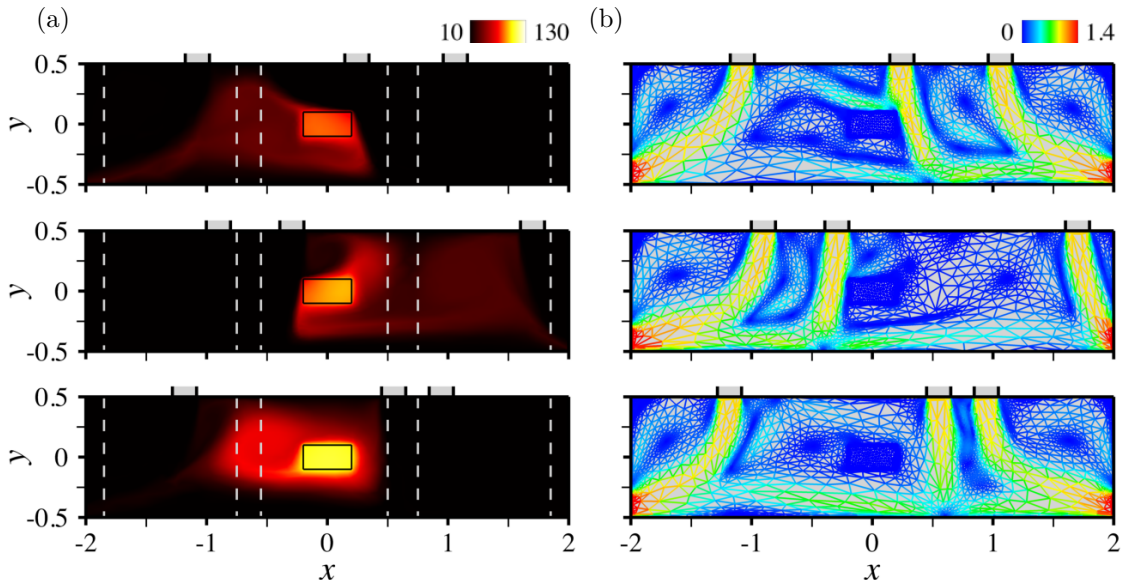


Figure 11: (a) Steady-state temperature against arrangements of 3 injectors, with admissible values under the fixed domain decomposition strategy S_1 delimited by the dashed lines. (b) Adapted meshes colored by the magnitude of velocity.

546 which especially yields $r_t = 0$ for a perfectly homogeneous cooling.

547 All results reported in the following are for $n_j = 3$ injectors. The agent is a fully-connected
 548 network with two hidden layers, each holding 2 neurons. The resolution process uses 8 environments
 549 and 2 steps mini-batches to update the network for 32 epochs, with learning rate set to 5×10^{-3} ,
 550 and PPO loss clipping range to $\epsilon = 0.3$.

551 5.3. Results

552 5.3.1. Fixed domain decomposition strategy

553 We consider first the so-called fixed domain decomposition strategy S_1 in which the top cavity
 554 wall is split into n_j equal subdomains, and each injector is forced to sit in a different subdomain (a
 555 somehow heavily constrained optimization problem if n_j is not too small, relevant for cases where
 556 the design is rigid and the practitioner has limited freedom to act). The edge values for the position
 557 x_k of the k -th injector read

$$x_k^- = -x_m + (k-1) \frac{2x_m + e_i}{n_j}, \quad x_k^+ = x_k^- + \frac{2x_m - (n_j - 1)e_i}{n_j}. \quad (41)$$

558 It can be checked that $x_k^- = x_{k-1}^+ + e_i$, so, it is possible to end up with two side-by-side injectors,
 559 which is numerically equivalent to having $n_j - 1$ injectors, $n_j - 2$ of width e_i plus one of width $2e_i$.
 560 For this case, 60 episodes have been run, each of which performs 1500 iterations with time step
 561 $\Delta t = 0.1$ to march the same initial condition (consisting of zero velocity and uniform temperature,
 562 except in the solid domain) to steady state, using the level set, velocity and temperature as multiple-
 563 component criterion to adapt the mesh (initially pre-adapted using the sole level set) every 5
 564 time steps under the constraint of a fixed number of elements $n_{el} = 15000$. This represents 480
 565 simulations, each of which is performed on 8 cores and lasts 10mn, hence 80h of total CPU cost.

566 It is out of the scope of this work to analyze in details the many flow patterns that develop when
 567 the blown fluid travels through the cavity. Suffice it to say that the outcome depends dramatically
 568 on the injectors arrangement, and features complex rebound phenomena (either fluid/solid, when
 569 a jet impinges on the cavity walls or on the workpiece itself, or fluid/fluid, when a deflected jet
 570 meets the crossflow of another jet), leading to the formation of multiple recirculation varying in
 571 number, position and size. Several such cases are illustrated in figure 11 via iso-contours of the
 572 steady-state temperature distributions, together with the corresponding adapted meshes colored by
 573 the magnitude of velocity to illustrate the ability of the numerical framework to capture accurately
 574 all boundary layers and shear regions via extremely stretched elements.

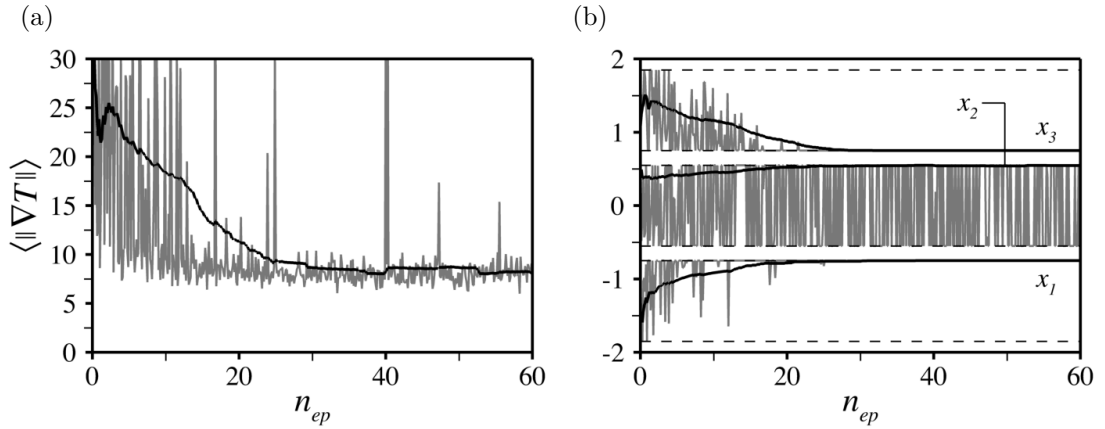


Figure 12: (a) Evolution per learning episode of the instant (in grey) and moving average (in black) rewards under the fixed domain decomposition strategy S_1 . (b) Same as (a) for the injectors center positions, with admissible values delimited by the dashed lines.

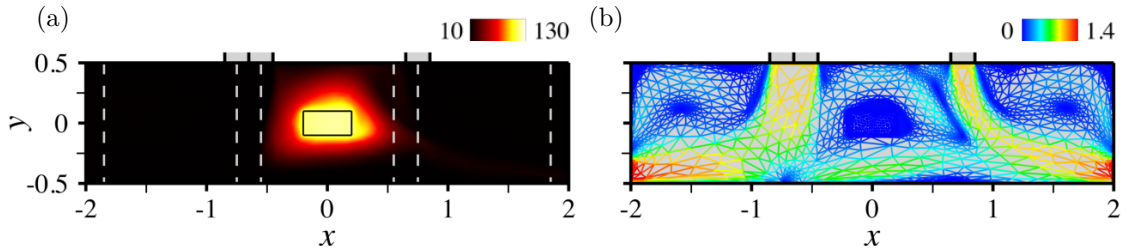


Figure 13: Same as figure 11 for the optimal arrangement of 3 injectors under the fixed domain decomposition strategy S_1 .

575 One point worth mentioning is that the individual position signals are best suited to draw robust
576 quantitative conclusion, as there is noise in the reward signal shown in figure 12(a). We believe
577 the issue to be twofold: on the one hand, the reward is approximated from point-wise temperature
578 data (similar to experimental measurements) that are more sensitive to small numerical errors
579 (e.g., the interpolation error at the probes position) than an integral quantity. On the other hand,
580 the mesh adaptation procedure is not a deterministic process, as the outcome depends on the
581 processors and number of processors used, and any initial difference propagates over the course
582 of the simulation because the meshes keep being adapted dynamically. In return, two exact same
583 control parameters can thus yield different rewards on behalf of different interpolation errors at
584 the probes position. This likely slows down learning and convergence, but we show in figure 12(b)
585 that the moving average distribution does converge to an optimal arrangement after roughly 25
586 episodes. The latter consists of an injector at the right-end of the left subdomain ($x_1^* = -0.75$)
587 and two side-by-side injectors sitting astride the center and right subdomains ($x_2^* = 0.55$ and
588 $x_3^* = 0.75$), that enclose the workpiece in a double-cell recirculation; see figure 13. These values
589 have been computed by averaging the instant positions of each injector over the 10 latest episodes,
590 with variations ± 0.002 computed from the root-mean-square of the moving average over the same
591 interval, a procedure that will be used consistently to assess convergence for all cases reported in the
592 following. The efficiency of the control itself is estimated by computing the magnitude of tangential
593 heat flux averaged over the same interval, found to be $\langle \|\nabla_{\parallel} T\| \rangle^* \sim 8.3$. Note, the position x_2^*
594 is actually obtained by averaging the absolute value of the instant position x_2 (although the true,
595 signed value is depicted in the figure), as the center injector keeps oscillating between two end
596 positions ± 0.55 on behalf of reflectional symmetry with respect to the vertical centerline.

597 5.3.2. Follow-up strategy

598 A less constrained problem is considered here using the so-called follow-up strategy S_2 , in which
599 all injectors are distributed sequentially the ones with respect to the others. The corresponding

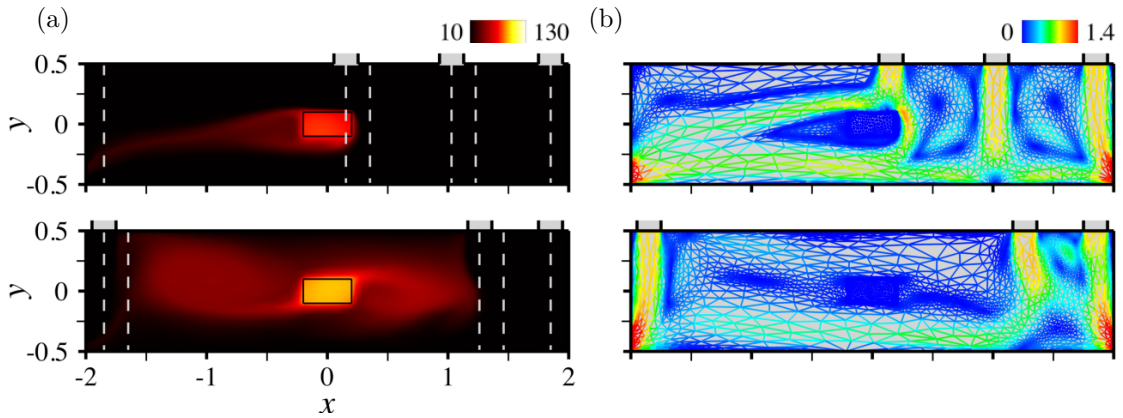


Figure 14: (a) Steady-state temperature against arrangements of 3 injectors, with admissible values under the follow-up strategy S_2 delimited by the dashed lines. (b) Adapted meshes colored by the magnitude of velocity.

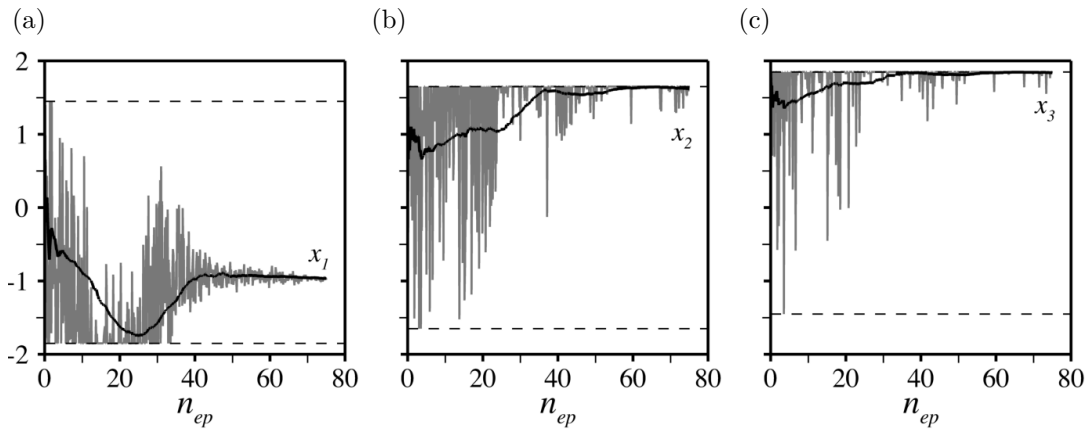


Figure 15: Evolution per learning episode of the instant (in grey) and moving average (in black) injectors center positions under the follow-up strategy S_2 , with admissible values delimited by the dashed lines.

600 edge values

$$x_1^- = -x_m, \quad x_1^+ = x_m - (n_j - 1)e_i, \quad (42)$$

$$x_k^- = x_{k-1}^+ + e_i, \quad x_k^+ = x_m - (n_j - k)e_i, \quad (43)$$

601 readily express that the k -th injector is forced to sit between the $k - 1$ -th one and the upper-right
602 cavity edge while leaving enough space to distribute the remaining $n_j - k$ injectors, which increases
603 the size of the control parameter space while again leaving the possibility for side-by-side injectors
604 (since $x_k^- = x_{k-1}^+ + e_i$ by construction). 75 episodes have been run for this case following the
605 exact same procedure as above, i.e., marching the zero-initial condition in time up to $t = 150$ with
606 $\Delta t = 0.1$, hence 600 simulations, each of which is performed on 8 cores and lasts 10mn, hence 100h
607 of total CPU cost.

608 The computed flow patterns closely resemble those obtained under the previous fixed domain
609 decomposition strategy, although figure 14 exhibits increased dissymmetry when two or more in-
610 jectors move simultaneously to the same side of the cavity. We show in figure 15 that the moving
611 average distribution converges after roughly 60 episodes, with the optimal arrangement consisting
612 of one injector roughly midway between the left cavity sidewall and the workpiece ($x_1^* = -0.96$),
613 and two side-by-side injectors at the right end of the cavity ($x_2^* = 1.65$ and $x_3^* = 1.85$). The
614 variations over the same interval are by ± 0.006 ; see also figure 16 for the corresponding flow pat-
615 tern. Convergence here is much slower than under S_1 , as the search for an optimal is complicated
616 by the fact that all injector positions are interdependent the ones on the others and it is up to the
617 network to figure out exactly how. Another contingent matter is that the agent initially spans a
618 fraction of the control parameter space because the large values of x_1 considered limit the space
619 available to distribute the other two injectors. This is all the more so as such configurations turn

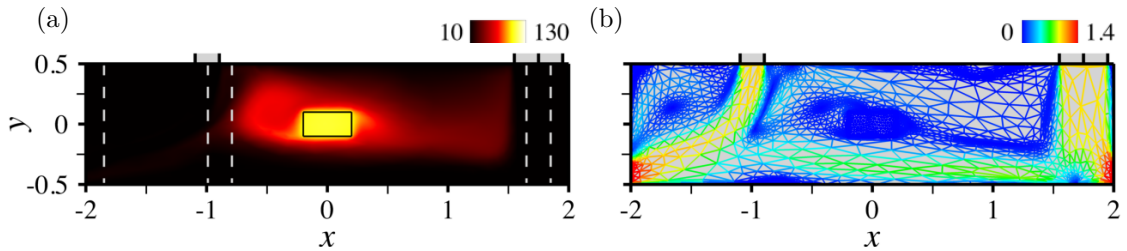


Figure 16: Same as figure 14 for the optimal arrangement of 3 injectors under the follow-up strategy S_2 .

620 to be far from optimality, for instance the magnitude of tangential heat flux is $\langle \|\nabla_{\parallel} T\| \rangle \sim 41.3$ for
 621 $x_1 = 1.45$, $x_2 = 1.65$ and $x_3 = 1.85$, but $\langle \|\nabla_{\parallel} T\| \rangle^* \sim 6.3$ at optimality. The latter value is smaller
 622 than the optimal achieved under S_1 , consistently with the fact that all positions spanned under S_1
 623 are admissible under S_2 , hence the S_1 optimal is expected to be a S_2 sub-optimal.

624 5.3.3. Free strategy

625 We examine now a third strategy S_3 referred to as the free strategy, in which all injectors are
 626 independent and free to move along the top cavity wall (a mildly constrained optimization problem,
 627 relevant for cases where the design is flexible and the practitioner has great freedom to act). The
 628 edge values for the position x_k of the k -th injector read

$$x_k^- = -x_m, \quad x_k^+ = x_m, \quad (44)$$

629 so two injectors can end up side-by-side and even overlapping one another if $|x_l - x_m| < e_i$. If
 630 so, we implement a single injector of width $e_i + |x_l - x_m|$ and maintain the blowing velocity (not
 631 the flow rate) for the purpose of automating the set-up design process, meaning that having n_j
 632 injectors, two of which overlap exactly (i.e., $|x_l - x_m| = 0$) is rigorously equivalent to having $n_j - 1$
 633 injectors. 60 episodes have been run for this case following the exact same procedure as above.

634 All flow patterns are reminiscent of those obtained under the previous fix decomposition S_1
 635 and follow-up S_2 strategies, even when two injectors overlap; see figure 17. Other than that, we
 636 show in figures 18 that the moving average distribution converges to an optimal consisting of two
 637 injectors almost perfectly overlapping one another at the left end of the cavity ($x_1^* = -1.85$ and
 638 $x_2^* = -1.82$), and a third injector at the right end of the cavity ($x_3^* = 1.85$). The variations over
 639 the same interval are by ± 0.007 , and the associated flow pattern shown in figure 19 is symmetrical
 640 and features two large recirculation regions on either side of the workpiece. Convergence occurs
 641 after roughly 40 episodes, i.e., faster than under S_2 (consistently with the fact that there is no
 642 need to learn anymore about how the network outputs depend the ones on the others) but slower
 643 than under S_1 (consistently with the fact that the size of the control parameter space has increased
 644 substantially). It is worth noticing that the system is invariant by permutations of the network
 645 outputs, meaning that there exist $2^{n_j} - 2$ distributions (hence 6 for $n_j = 3$) associated with the
 646 same reward. Nonetheless, a single optimal is selected, which is essentially fortuitous since the
 647 agent does not learn about symmetries under the optimization process (otherwise S_1 would have
 648 similarly selected a single optimal). The magnitude of tangential heat flux is $\langle \|\nabla_{\parallel} T\| \rangle^* \sim 11.2$ at
 649 optimality, i.e., larger than that achieved under S_2 . This can seem surprising at first, because all
 650 positions spanned under S_2 are admissible under S_3 , and the S_2 optimal is thus expected to be a
 651 S_3 sub-optimal. However, the argument does not hold here because the overlap in the S_3 optimal
 652 reduces the flow rate to that of a two-injectors set-up, so the comparison should be with the S_2
 653 optimal with $n_j = 2$.

654 5.3.4. Inverse strategy

655 Finally, we propose here to make the most of the numerical framework flexibility to solve a
 656 different optimization problem consisting in selecting first an injector distribution, then in finding
 657 the position x_0 of the solid center of mass minimizing the magnitude of tangential heat flux (which
 658 is relevant for cases where the practitioner simply cannot act on the design). The so-called inverse
 659 strategy S_4 considered herein features two injectors at each end of the cavity ($x_1 = -1.85$ and
 660 $x_2 = 1.85$), identical to the optimal arrangement of 3 injectors under the free strategy S_3 . The

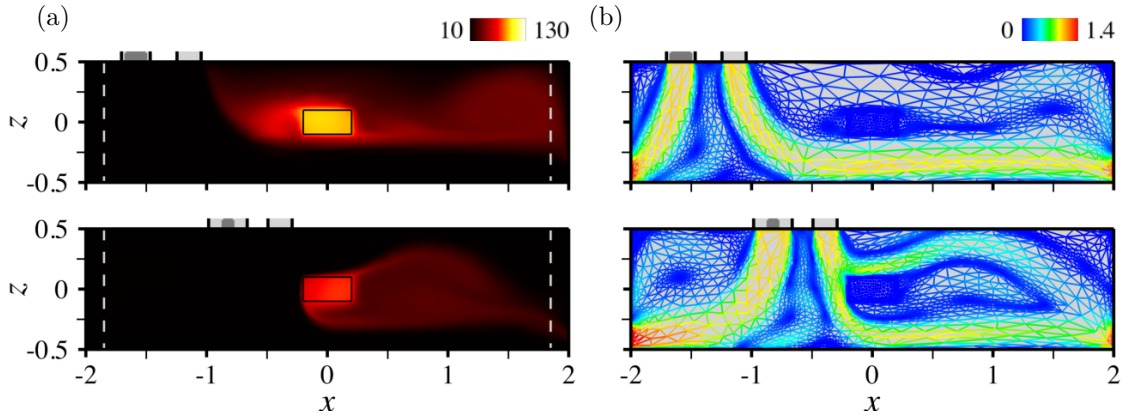


Figure 17: (a) Steady-state temperature against arrangements of 3 injectors, with admissible values under the free strategy S_3 delimited by the dashed lines and overlaps marked by the dark grey shade. (b) Adapted meshes colored by the magnitude of velocity.

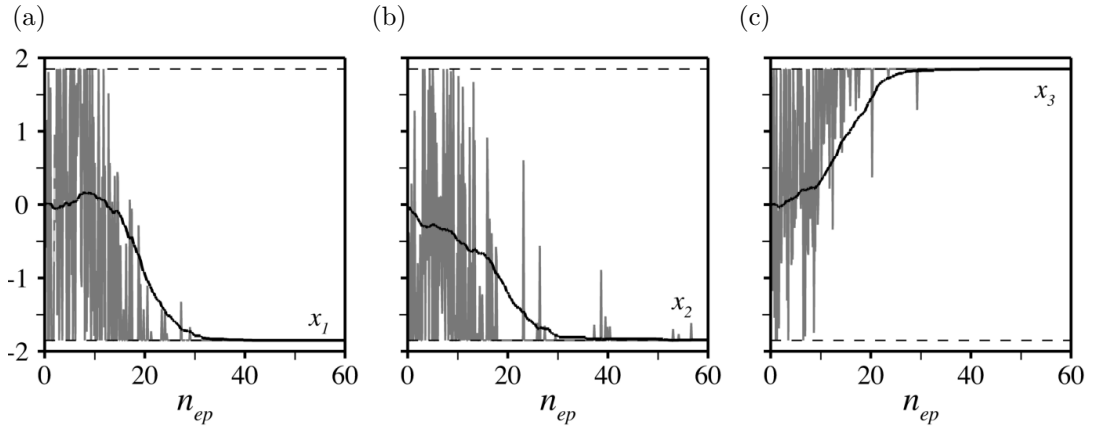


Figure 18: Evolution per learning episode of the instant (in grey) and moving average (in black) injectors center positions under the free strategy S_3 , with admissible values delimited by the dashed lines.

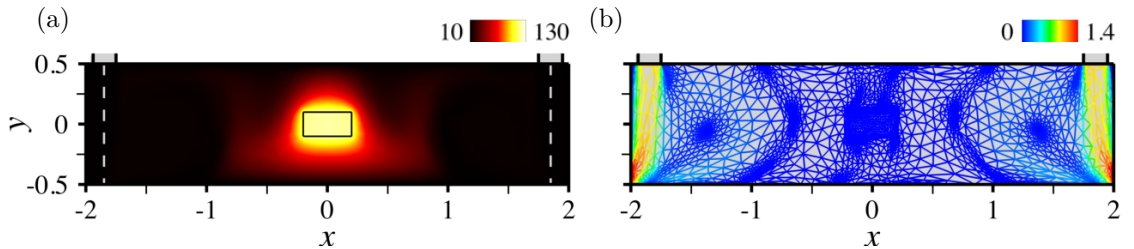


Figure 19: Same as figure 17 for the optimal arrangement of 3 injectors under the free strategy S_3 .

661 center of mass can take any value in $[-x_{0m}; x_{0m}]$ where we set $x_{0m} = 2(H - h)$ to avoid numerical
 662 issues at the sidewalls. The same coordinate system as above is used, but with reference frame
 663 attached to the cavity, not the moving solid (hence all results obtained under the previous strategies
 664 pertain to $x_0 = 0$ in the new system).

665 A total of 60 episodes have been run for this case using the exact same DRL agent, the only
 666 difference being in the network action output, now made up of a single value $\hat{x}_0 \in [-1; 1]$, mapped
 667 into the actual position using

$$x_0 = x_{0m} \hat{x}_0. \quad (45)$$

668 A large variety of flow patterns is obtained by doing so, that closely resemble those computed under
 669 the previous strategies, only the outcome is now also altered by the width of the gap between the
 670 cavity sidewalls and the workpiece, as illustrated in figure 20. We show in figure 21 that the position
 671 of the solid center of mass converges to an optimal $x_0^* = 0.42$ (the variations over the same interval

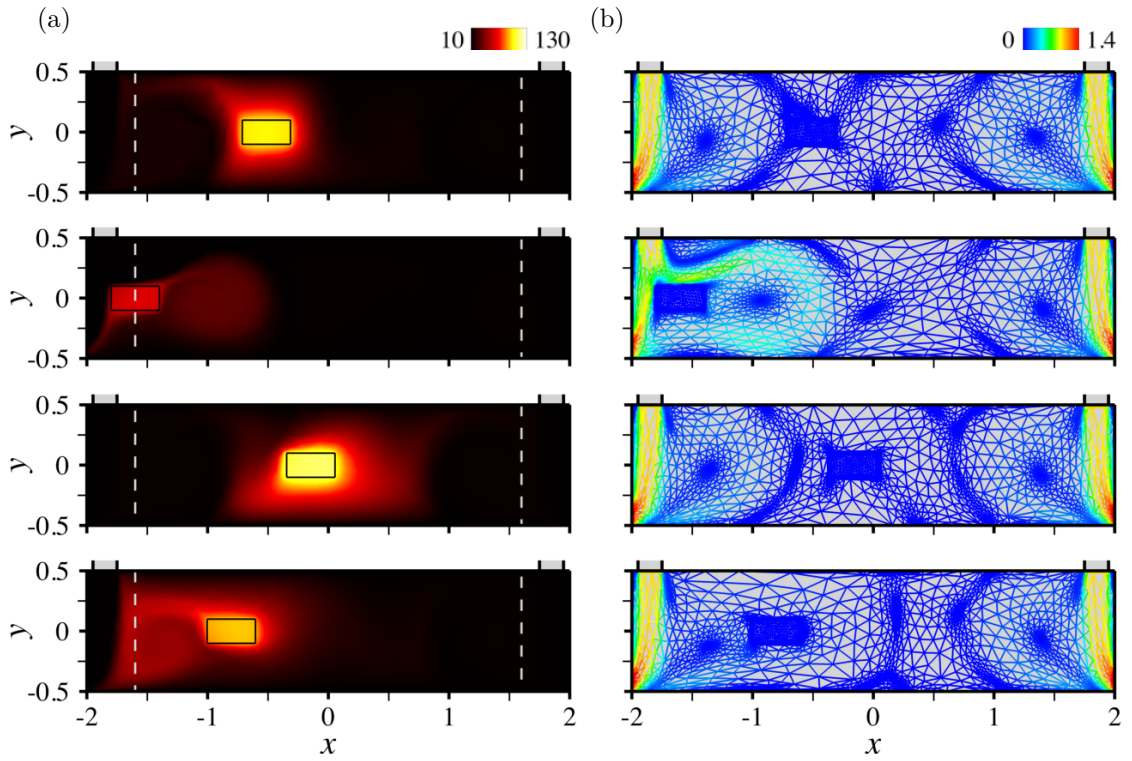


Figure 20: (a) Steady-state temperature against solid center of mass position, with admissible domains under the inverse strategy S_4 marked by the dashed lines. (b) Adapted meshes colored by the norm of velocity.

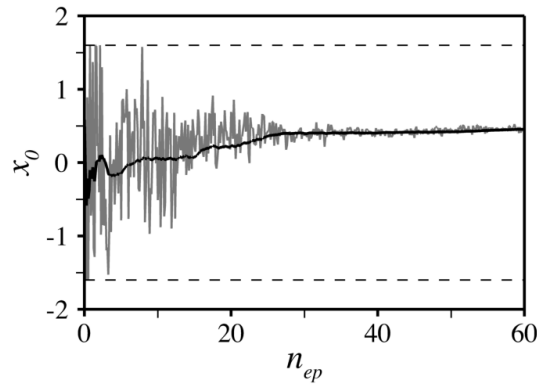


Figure 21: Evolution per learning episode of the instant (in grey) and moving average (in black) center of mass positions under the inverse strategy S_4 , with admissible values delimited by the dashed lines.

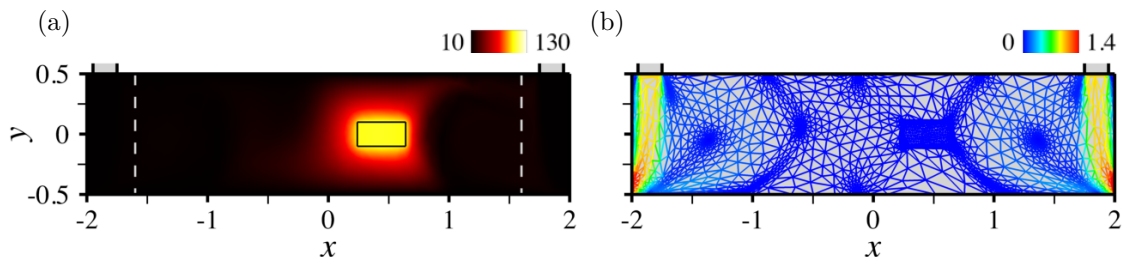


Figure 22: Same as figure 20 for the optimal center of mass position under the inverse strategy S_4 .

672 being by ± 0.005), the associated magnitude of tangential heat flux $\langle \|\nabla_{\parallel} T\| \rangle^* \sim 4.1$, being smaller
 673 than that achieved under S_3 using a centered workpiece. The fact that the optimal position is
 674 offset from the vertical centerline is a little surprising at first, because intuition suggests that the

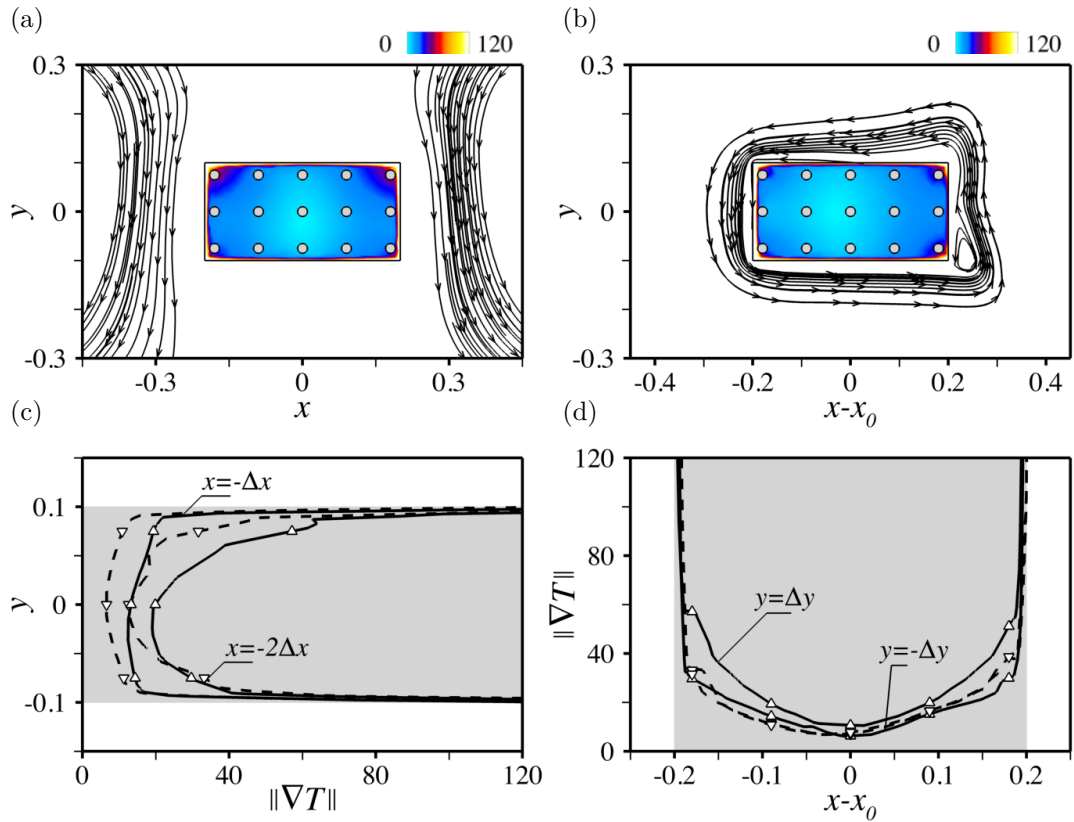


Figure 23: (a,b) Norm of the temperature gradient in the solid domain with superimposed streamlines of the underlying velocity field, as computed for (a) $x_0 = 0$, and (b) $x_0^* = 0.45$, i.e., the optimal position selected under the inverse strategy S_4 . (c) Cuts along the two leftmost columns of probes. The solid and dashed lines refer to $x_0 = 0$ and $x_0^* = 0.42$, respectively, and the symbols mark the probe values. (d) Same as (c) for cuts along the lower and upper rows of columns.

675 simplest way to achieve homogeneous heat transfer is by having symmetrically distributed injectors.
676 Nonetheless, examining carefully the norm of the temperature gradient in the solid domain shows
677 that $x_0 = 0$ achieves close to perfect horizontal symmetry but vertical asymmetry, owing to the
678 formation of two large-scale, small velocity end vortices entraining heat laterally downwards; see
679 figure 23(a). Conversely, for $x \sim x_0^*$, the workpiece is almost at the core of the closest recirculation
680 region, hence the surrounding fluid particles have small velocities and wrap almost perfectly around
681 its surface, as illustrated in figure 23(b). This restores excellent vertical symmetry, as evidenced
682 by relevant cuts along the two leftmost columns of probes in figure 23(c), and along the lower and
683 upper rows in figure 23(d), which explains the improved reward.

684 5.4. Discussion

685 Figure 24 reproduces the optimal temperature distributions computed under the various strate-
686 gies considered above. For benchmarking purposes, we also provide in table 3 relevant convergence
687 data computed over the 10 latest episodes. To recap, the most homogeneous cooling is achieved
688 under the follow-up strategy S_2 , but the DRL agent seems more easily trained under the fixed
689 decomposition domain strategy S_1 and the free strategy S_3 . Another interesting point is the ex-
690 tent to which the workpiece is actually cooled, for which S_2 seems more relevant, on behalf of the
691 dissymmetry in the left and right flow rates that creates order one velocities at the bottom of the
692 cavity. This stresses S_2 as a possible compromise to achieve efficient *and* homogeneous cooling,
693 although a true optimal with this regard can be computed rigorously by applying the same ap-
694 proach to compound functionals weighing, e.g., the magnitude of the tangential heat flux and the
695 solid center temperature (which we defer to future work).

696 These results provide a basis for future self-assessment of the method and identifies potential
697 for improvement regarding the convergence efficiency. The approach can certainly benefit from

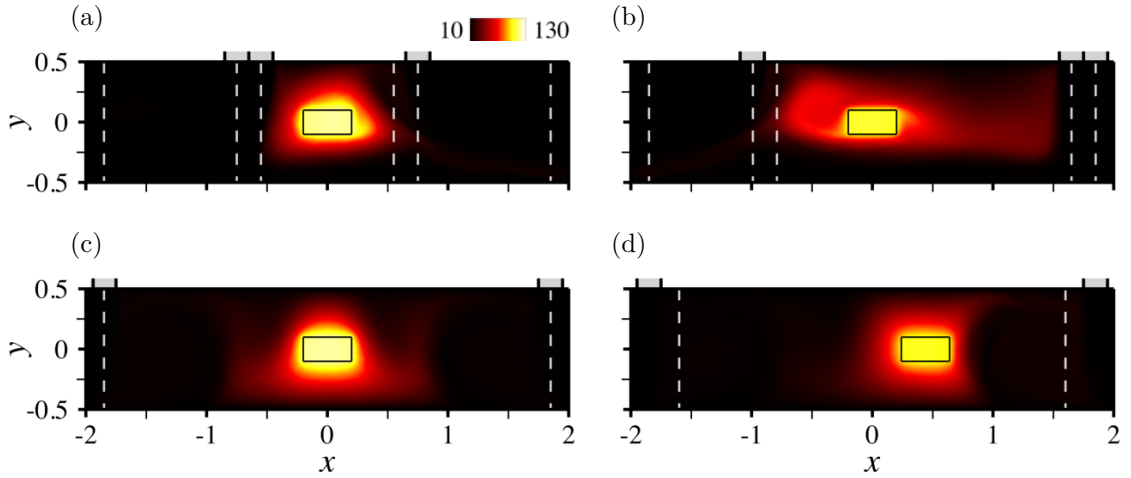


Figure 24: (a-c) Optimal arrangements of 3 injectors under the (a) fixed decomposition domain strategy S_1 , (b) follow-up strategy S_2 and (c) free strategy S_3 . (d) Optimal position of the workpiece under the inverse strategy S_4 .

	n_j	n_{ep}	x_0	x_1	x_2	x_3	$\langle \ \nabla_{\parallel} T\ \rangle$
S_1	3	60	0	-0.75	± 0.55	0.75	8.3
S_2	3	75	0	-0.96	1.65	1.85	6.3
S_3	3	60	0	-1.85	-1.82	1.85	11.2
S_4	2	60	0.42	-1.85	1.85	-	4.1

Table 3: Numerical data for the optimal arrangements computed under strategies S_1 -4. All values computed by averaging the instant signal over the 10 latest learning episodes.

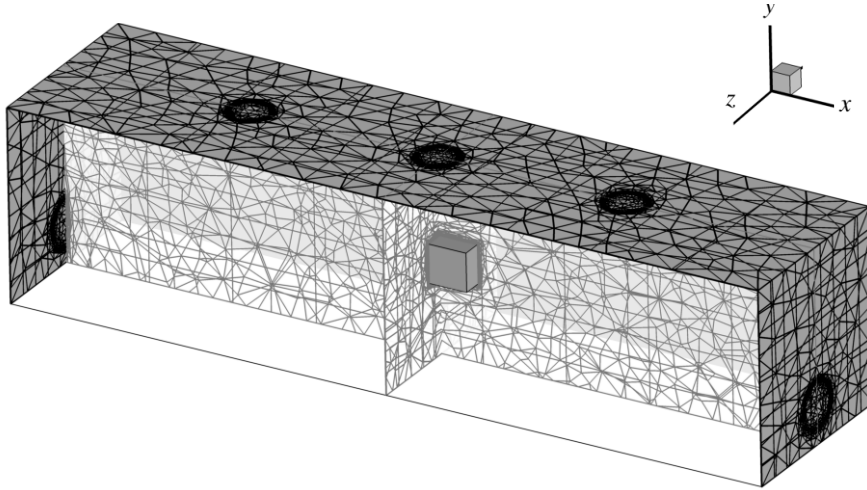


Figure 25: Schematic of the 3-D forced convection set-up.

698 a fine tuning of the reward computation, as having sufficient spatial resolution on the relevant
699 state of the system is an obvious requirement to allow a successful control. Adjusting the trade-off
700 between exploration and exploitation is also worth consideration to better handle the existence of
701 multiple global optima (whether they stem from symmetries or from the topology of the reward
702 itself) which could be done using non-normal probability density functions.

H	h	d_i	d_0	δ_0	V_i	T_w	T_c	T_h	μ	ρ	λ	c_p	
1	0.2	0.2	0.24	0.16	1	10	10	150	0.01	1	0.5	1000	Fluid
									1000	100	15	300	Solid

Table 4: Numerical parameters used in the 3-D forced convection problem. All values in SI units, with the exception of temperatures given in Celsius.

6. Extension to 3-D forced convection

6.1. Case description

The model cooling set up considered in section 5 is extended here to 3-D to assess the extent to which the approach carries over to three-dimensional conjugate heat transfer. The main differences between 2-D and 3-D are as follows: a Cartesian coordinate system is used with origin at the center of mass of the solid, horizontal x -axis, vertical y -axis, and the z -axis completes the direct triad; see figure 25. The solid is a rectangular prism with aspect ratio 2:1:1, and is fixed at the center of a rectangular cavity with height H and aspect ratio 4:1:1. We consider n_j circular-shaped injectors with diameter d_i , whose exit planes are forced to be symmetrical with respect to $z = 0$, hence each injector is identified by the horizontal position of its center $x_{k \in \{1 \dots n_j\}}$. We also use circular-shaped exhaust areas with diameter d_o , offset by a distance δ_o from the bottom of the cavity, and whose exit planes are also symmetrical with respect to $z = 0$, hence each exhaust area is identified by the vertical position of its center $(d_o + \delta_o - H)/2$. The governing equations are solved with the exact same boundary conditions as in section 5. All parameters are provided in Table 4, including the material properties used to model the composite fluid, that yield fluid values of the Reynolds and Prandtl numbers

$$\text{Re} = \frac{\rho V_i d_i}{\mu} = 20, \quad \text{Pr} = 20. \quad (46)$$

6.2. Control strategy

We keep here the same control objective and compute the reward fed to the DRL from 45 probes arranged symmetrically into $n_z = 3$ transverse layers with resolution $\Delta z = 0.075$, each of which distributes uniformly 15 probes into $n_x = 5$ columns and $n_y = 3$ rows with resolutions $\Delta x = 0.09$ and $\Delta y = 0.075$. In practice, the 3-D reward is simply the average over z of the 2-D reward defined in section 5, hence $r_t = -\langle \|\nabla_{\parallel} T\| \rangle$ with

$$\langle \|\nabla_{\parallel} T\| \rangle = \frac{1}{(n_x + n_y)n_z} \sum_{i,j,k} \langle \|\nabla_{\parallel} T\| \rangle_{ik} + \langle \|\nabla_{\parallel} T\| \rangle_{jk}, \quad (47)$$

with

$$\langle \|\nabla_{\parallel} T\| \rangle_{ik} = \frac{2}{n_y - 1} \left| \sum_{j \neq 0} \text{sgn}(j) \|\nabla T\|_{ijk} \right|, \quad \langle \|\nabla_{\parallel} T\| \rangle_{jk} = \frac{2}{n_x - 1} \left| \sum_{i \neq 0} \text{sgn}(i) \|\nabla T\|_{ijk} \right|, \quad (48)$$

and the subscripts ik , jk and ijk denote quantities evaluated at $(x, z) = (i\Delta x, k\Delta z)$, $(y, z) = (j\Delta y, k\Delta z)$ and $(x, y, z) = (i\Delta x, j\Delta y, k\Delta z)$, respectively.

All results reported in the following are for $n_j = 3$ injectors. The edge values needed to map the network action output into the actual injectors positions deduce straightforwardly from (41)-(44) substituting the diameter d_i of the 3-D injectors for the length e_i of the 2-D injectors. The same DRL agent is used, that consists of two hidden layers, each holding 2 neurons, and the resolution process uses 8 environments and 2 steps mini-batches to update the network for 32 epochs. Each environment performs 1250 iterations with time step $\Delta t = 0.1$ to march the same initial condition (consisting of zero velocity and uniform temperature, except in the solid domain) to steady state, using the level set, velocity and temperature as multiple-component criterion to adapt the mesh (initially pre-adapted using the sole level set) every 10 time steps under the constraint of a fixed

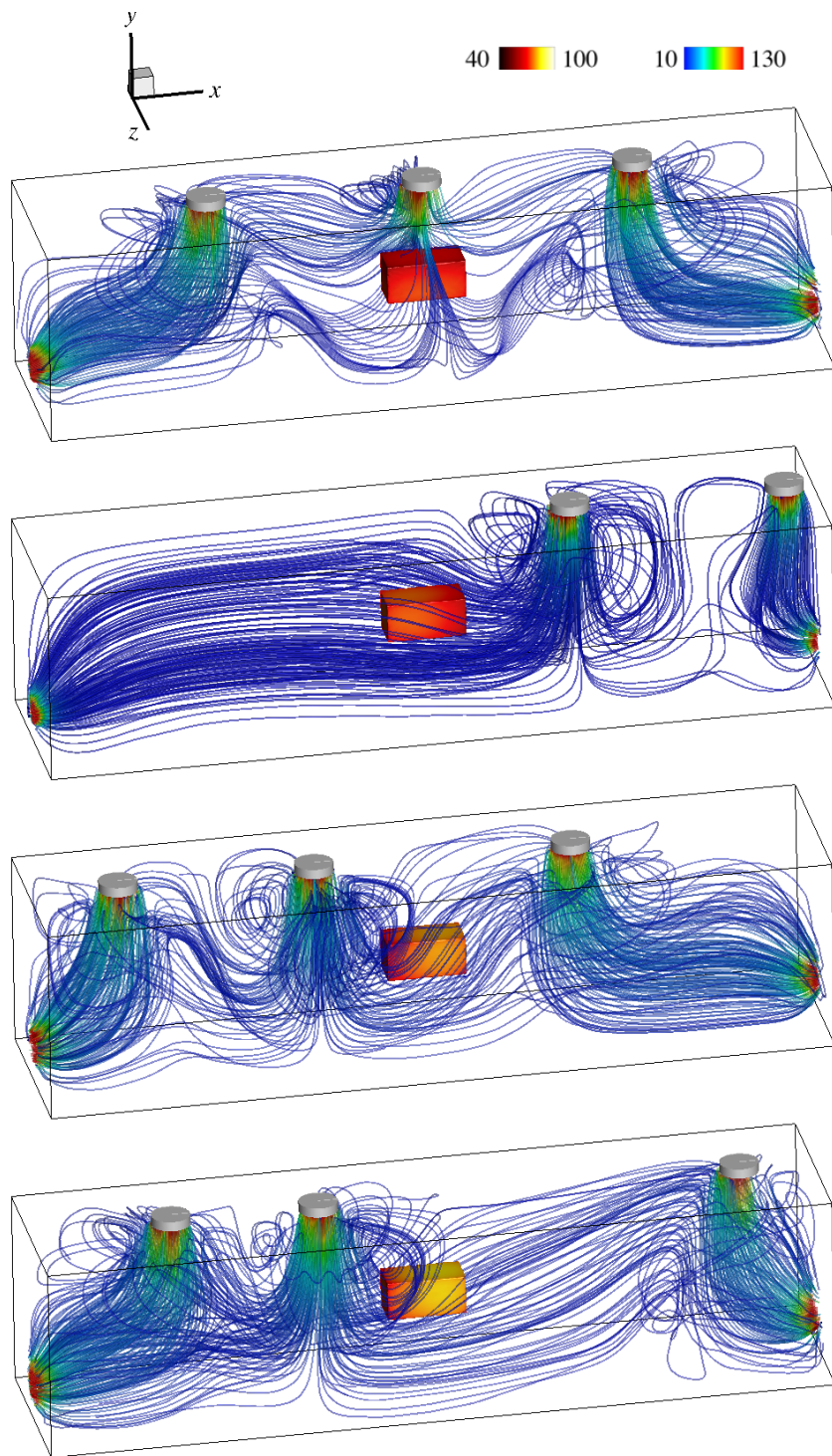


Figure 26: Representative steady-state temperature distributions at the solid/fluid interface together with 3-D streamlines colored by the magnitude of velocity.

737 number of elements $n_{el} = 120000$. This is likely insufficient to claim true numerical accuracy, but
 738 given the numerical cost (320 3-D simulations per strategy, each of which is performed on 8 cores
 739 and lasts 2h30, hence 800h of total CPU cost), we believe this is a reasonable compromise to assess
 740 feasibility while producing qualitative results to build on.

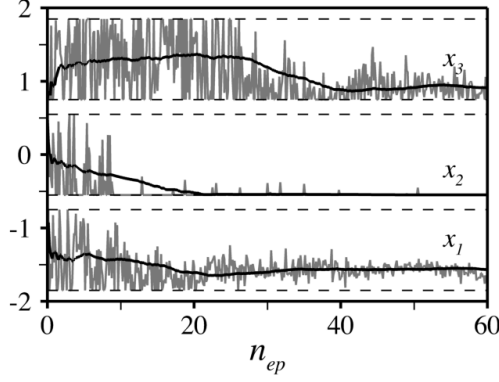


Figure 27: Evolution per learning episode of the instant (in grey) and moving average (in black) injectors center positions under the three-dimensional fixed domain decomposition strategy S_1 , with admissible values delimited by the dashed lines.

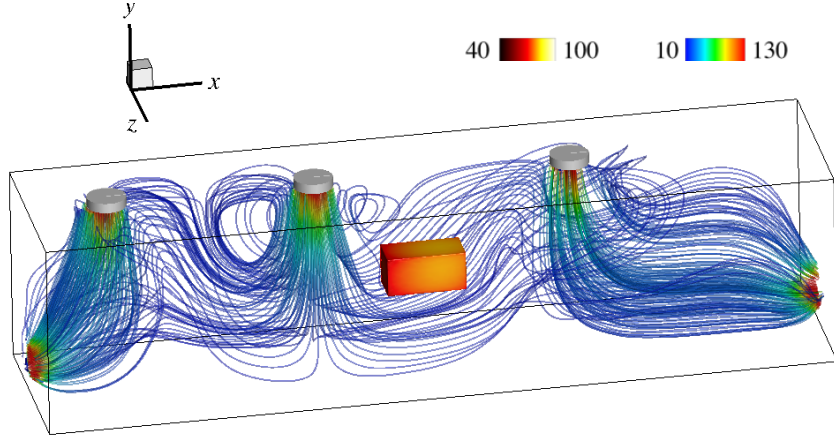


Figure 28: Optimal 3 injector arrangement under the three-dimensional fixed decomposition domain strategy S_1 .

741 6.3. Results

742 Only the fixed domain decomposition S_1 strategy (in which the top cavity wall is split into
 743 n_j equal subdomains and each injector is forced to sit in a different subdomain) and the free S_3
 744 strategy (in which the injectors are entirely independent and free to move along the top cavity
 745 wall) are considered here to save computational resources, as learning has been seen to be slower
 746 in 2-D under the follow-up S_2 strategy.

747 A total of 60 episodes have been run under the fixed domain decomposition strategy S_1 . Several
 748 representative flow patterns computed over the course of optimization are shown in figure 26 via
 749 iso-contours of the steady-state temperature at the fluid-solid interface and 3-D streamlines colored
 750 by the magnitude of velocity, to put special emphasis on transverse inhomogeneities and display
 751 the increased degree of complexity due to the formation of large-scale horseshoe vortices wrapped
 752 around the nozzle jets. We show in figure 27 that the distribution slowly converges to an optimal
 753 arrangement consisting of one injector at the left end of the left subdomain ($x_1^* = -1.63$), another
 754 one at the left end of the center subdomain ($x_2^* = -0.55$), and a third one at the left end of
 755 the right subdomain ($x_3^* = 0.87$), as has been determined by averaging the instant positions
 756 of each injector over the latest 10 learning episodes, with variations by roughly ± 0.04 computed
 757 from the root-mean-square of the moving average over the same interval. This is larger by one
 758 order of magnitude than the variations reported in 2-D, as the agent keeps exploring slightly
 759 sub-optimal positions of the lateral injectors, which likely simply reflects the challenging nature
 760 of performing three-dimensional optimal control. The 3-D S_1 optimal somehow resemble its 2-D
 761 counterpart, namely the center injector is at the exact same position, while the lateral injectors
 762 (especially the leftmost one) have been pushed towards the cavity sidewalls. The associated flow
 763 pattern is reported in figure 28. The associated optimal reward computed over the same interval

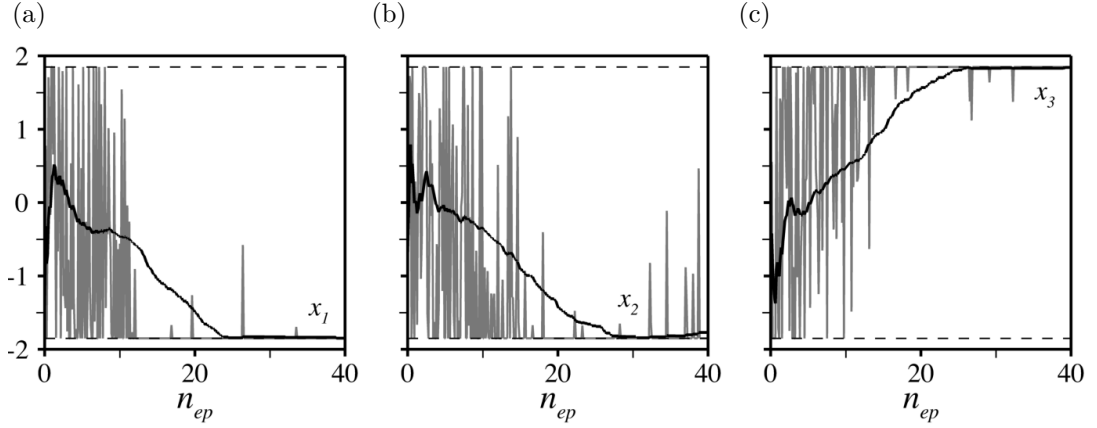


Figure 29: Evolution per learning episode of the instant (in grey) and moving average (in black) injectors center positions under the three-dimensional free strategy S_3 , with admissible values delimited by the dashed lines.

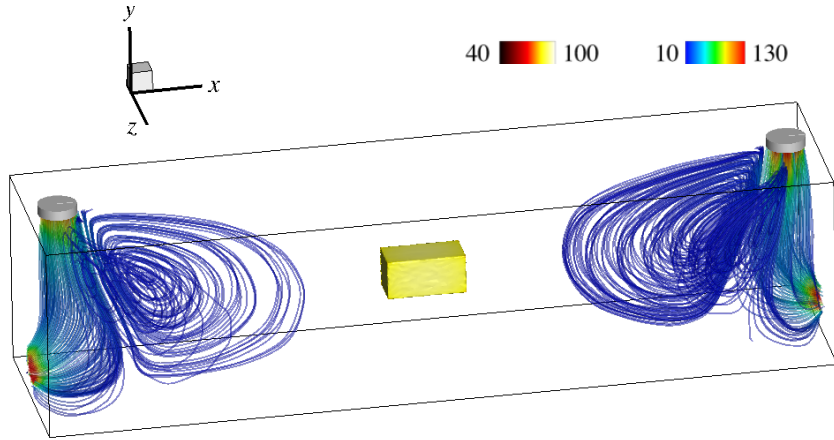


Figure 30: Optimal 3 injector arrangement under the three-dimensional free strategy S_3 .

	n_j	n_{ep}	x_0	x_1	x_2	x_3	$\langle \ \nabla_{\parallel} T\ \rangle$
S_1	3	60	0	-1.63	-0.55	0.87	19.5
S_3	3	40	0	-1.83	-1.82	1.83	4.7

Table 5: Numerical data for the optimal arrangements computed in three-dimensions under strategies S_1 and S_3 . All values computed by averaging the instant signal over the 10 latest learning episodes.

764 is $\langle \|\nabla_{\parallel} T\| \rangle^* \sim 19.5$, i.e. twice as large than in 2-D, although it is difficult to compare further
 765 because of the difference in the Reynolds and Prandtl number.

766 Another 40 episodes have been run under the free strategy S_3 , for which the results are al-
 767 most identical to their 2-D counterparts, as the distribution converges in figure 29 to an optimal
 768 arrangement consisting of two overlapping injectors at the left end of the cavity ($x_1^* = -1.83$ and
 769 $x_2^* = -1.82$), and a third injector at the right end ($x_3^* = 1.83$), with variations by with ± 0.01
 770 for the lateral injectors, but ± 0.03 for the center injector, for which the agent keeps occasionally
 771 exploring sub-optimal positions. The corresponding flow pattern shown in figure 30 is thus again
 772 symmetrical with two large, 3-D recirculation regions on either side of the workpiece. The asso-
 773 ciated optimal reward computed over the same interval is $\langle \|\nabla_{\parallel} T\| \rangle^* \sim 4.7$ substantially smaller
 774 than that achieved under the 3-D S_1 strategy, which again demonstrates the feasibility to improve
 775 performances by allowing overlaps. All relevant numerical data are reported in table 5 for the sake
 776 of completeness.

777 **7. Conclusion**

778 Optimization of conjugate natural and forced heat transfer systems is achieved here training a
779 fully connected network with a novel single-step PPO deep reinforcement algorithm, in which it
780 gets only one attempt per learning episode at finding the optimal. The numerical reward fed to the
781 network is computed with a finite elements CFD environment solving stabilized weak forms of the
782 coupled Navier–Stokes and heat equations with a combination of variational multi-scale modeling,
783 immerse volume method, and multi-component anisotropic mesh adaptation.

784 Convergence is assessed by alleviating the natural convection induced enhancement of heat
785 transfer in a two-dimensional, differentially heated square cavity controlled by piece-wise constant
786 fluctuations of the sidewall temperature. The approach is also relevant to forced convection prob-
787 lems, as single-step PPO shows capable of improving the homogeneity of temperature across the
788 surface of two and three-dimensional hot workpieces under impingement cooling. Several control
789 strategies are considered, in which the position of multiple cold air injectors is optimized relative
790 to a fixed workpiece position, each of which mimics a different levels of design constraint. The
791 flexibility of the numerical framework also allows solving the inverse problem, i.e., optimizing the
792 workpiece position relative to a fixed injector distribution, which is relevant in situations where the
793 design cannot be changed. The approach is beneficial in two important respects: first, it is effi-
794 cient, even though the parameter spaces are large and it may be costly to identify optimal control
795 parameters from simple parametric searches. Second, and more significantly, it is capable of deter-
796 mining additional optimal configurations, as the results of the inverse problem under symmetrical
797 actuation indicate that the workpiece is best positioned offset from the symmetry axis, which had
798 not been anticipated. Such results clearly stress that single-step PPO (and DRL in general) can
799 be effective to explore and discover new solutions from unforeseen parameter combinations.

800 Fluid dynamicists have just begun to gauge the ability of DRL to design optimal control strate-
801 gies. The efforts for developing single-step PPO are ongoing and remain at an early stage, so we
802 do not expect the approach to compete right away with more established methods, for instance
803 Evolution strategies (ES), a popular class of algorithms imitating principles of organic evolution
804 processes as rules for black-box optimum seeking. ES rely on a stochastic description of the vari-
805 ables to optimize, i.e., they consider probability density functions, not deterministic variables.
806 Simply put, at each generation (or iteration) new candidate solutions are sampled isotropically
807 by variation of the current parental individuals according to a multivariate normal distribution.
808 Recombination and mutation transformations are applied (that amount respectively to changing
809 the mean and adding a random, zero-mean perturbation), after which the individuals with the
810 highest cost function are selected to become the parents in the next generation. Improved variants
811 include the covariance matrix adaptation evolution strategy (CMA-ES), that speeds up conver-
812 gence by updating its full covariance matrix (which amounts to learning a second-order model of
813 the objective function). In present form, single-step PPO can be thought as an evolutionary-like
814 algorithm with simpler heuristics (i.e., without an evolutionary update strategy, as the optimal
815 model parameters are learnt via gradient ascent), so it is our guess that the performance should
816 be comparable to that of standard ES algorithms with isotropic covariance matrix. Besides con-
817 solidating the acquired knowledge, future research should thus aim at improving efficiency (by
818 fine-tuning the hyper parameters, or using pre-trained deep learning models) and convergence (by
819 coupling with a surrogate model trained on-the-fly, using non-normal probability density functions,
820 or modifying the balance between exploration and exploitation, as PPO prevents large updates of
821 the policy to avoid the issue of performance collapse). **For complex configurations representative of
822 industrial applications, the implementation of properly designed numerical rewards (under partial
823 state information) and noise reduction techniques is another issue that deserves consideration, as
824 pointed out in [38].**

825 Scope is another key ingredient to keep pushing forward the state of the art. The next step
826 is to tackle more complex test cases exhibiting flow unsteadiness and turbulence, which the CFD
827 environment is perfectly suited to do via a combination of Reynolds-averaged Navier–Stokes mod-
828 eling [76, 77] and second-order, semi-implicit time discretization [78]. We believe that this will
829 highlight even more clearly the relevance of the methodology, as [42] speculates that DRL should
830 be able to handle chaotic systems without suffering from the shortcomings and limitations of the
831 adjoint method, and it is shown in [39] to outperform a canonical linear proportional-derivative

832 controller in controlling turbulent natural convection. The long-term objective would be to en-
833 rich the description of the test cases using multi-physics modeling (e.g., radiative heat transfer,
834 phase transformation) in order to pave the way toward flexible, ready-to-use control of industrially
835 relevant applications, such as thermal comfort for building design or manufacturing processes.

836 Acknowledgements

837 This work is supported by the Carnot M.I.N.E.S. Institute through the M.I.N.D.S. project.

838 References

- 839 [1] A. Jameson, Aerodynamic design via control theory, *J. Sci. Comput.* 3 (1998) 233–260.
- 840 [2] M. D. Gunzburger, *Perspectives in flow control and optimization*, SIAM, Philadelphia, 2002.
- 841 [3] T. Bewley, Flow control : new challenges for a new renaissance, *Prog. Aerosp. Sci.* 37 (2001)
842 21–58.
- 843 [4] K. Momose, K. Abe, H. Kimoto, Reverse computation of forced convection heat transfer for
844 optimal control of thermal boundary conditions, *Heat Tran. Asian Res.* 33 (2004) 161–174.
- 845 [5] A. Belmiloudi, Robin-type boundary control problems for the nonlinear Boussinesq type equa-
846 tions, *J. Math. Anal. Appl.* 273 (2002) 428–456.
- 847 [6] G. Bärwolff, M. Hinze, Optimization of semiconductor melts, *Z. Angew. Math. Mech.* 86
848 (2006) 423–437.
- 849 [7] J. Boldrini, E. Fernández-Cara, M. Rojas-Medar, An optimal control problem for a generalized
850 Boussinesq model: The time dependent case, *Rev. Mat. Comput.* 20 (2007) 339–366.
- 851 [8] H. Karkaba, T. Dbouk, C. Habchi, S. Russeil, T. Lemenand, D. Bougeard, Multi objective
852 optimization of vortex generators for heat transfer enhancement using large design space
853 exploration, *Chem. Eng. Process.* (2020 (accepted)).
- 854 [9] P. Meliga, J.-M. Chomaz, Global modes in a confined impinging jet: application to heat
855 transfer and control, *Theor. Comput. Fluid Dyn.* 25 (1) (2011) 179–193.
- 856 [10] D. Rumelhart, G. Hinton, R. Williams, Learning representations by back-propagating errors,
857 *Nature* 323 (1986) 533–536.
- 858 [11] J. Kober, J. A. Bagnell, J. Peters, Reinforcement learning in robotics: A survey, *The Interna-
859 tional Journal of Robotics Research* 32 (11) (2013) 1238–1274.
- 860 [12] V. Mnih, K. Kavukcuoglu, D. Silver, R. A.A., V. J., M. Bellemare, A. Graves, M. Riedmiller,
861 A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King,
862 D. Kumaran, D. Wierstra, S. Legg, D. Hassabis, Human-level control through deep reinforce-
863 ment learning, *Nature* 518 (2015) 7540.
- 864 [13] G. E. Hinton, A. Krizhevsky, I. Sutskever, Imagenet classification with deep convolutional
865 neural networks, *Advances in neural information processing systems* 25 (2012) 1106–1114.
- 866 [14] B. Lusch, J. N. Kutz, S. L. Brunton, Deep learning for universal linear embeddings of nonlinear
867 dynamics, *Nature communications* 9 (1) (2018) 1–10.
- 868 [15] M. Raissi, A. Yazdani, G. E. Karniadakis, Hidden Fluid Mechanics: A Navier-Stokes Informed
869 Deep Learning Framework for Assimilating Flow Visualization Data, *arXiv* (2018).
870 URL <http://arxiv.org/abs/1808.04327>
- 871 [16] A. D. Beck, D. G. Flad, C.-D. Munz, Deep Neural Networks for Data-Driven Turbulence
872 Models, *arXiv* (2018).
873 URL <http://arxiv.org/abs/1806.04482>

- 874 [17] S. L. Brunton, B. Noack, P. Koumoutsakos, Machine learning for fluid mechanics, *Annu. Rev.*
875 *Fluid Mech.* 52 (2020) 477–508.
- 876 [18] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert,
877 L. Baker, M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. van den Driess-
878 che, T. Graepel, D. Hassabis, Mastering the game of go without human knowledge, *Nature*
879 550 (7676) (2017) 354–359.
- 880 [19] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, M. Hutter, Learning
881 agile and dynamic motor skills for legged robots, *Science Robotics* 4 (26) (2019) eaau5872.
- 882 [20] A. Bernstein, E. Burnaev, Reinforcement learning in computer vision, *Proc. SPIE 10696*, 10th
883 International Conference on Machine Vision (ICMV 2017) (2018).
- 884 [21] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, D. Wierstra,
885 Continuous control with deep reinforcement learning, arXiv e-prints (2015).
886 URL <http://arxiv.org/abs/1509.02971>
- 887 [22] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, O. Klimov, Proximal Policy Optimization
888 Algorithms, arXiv e-prints (Jul. 2017). [arXiv:1707.06347](https://arxiv.org/abs/1707.06347).
- 889 [23] V. Belus, J. Rabault, J. Viquerat, Z. Che, E. Hachem, U. Réglade, Exploiting locality
890 and translational invariance to design effective deep reinforcement learning control of the
891 1-dimensional unstable falling liquid film, *AIP Adv.* 9 (2019) 125014.
- 892 [24] M. A. Bucci, O. Semeraro, A. Allauzen, G. Wisniewski, L. Cordier, L. Mathelin, Control of
893 chaotic systems by deep reinforcement learning, *Proc. Roy. Soc. A* 475 (2019) 20190351.
- 894 [25] G. Novati, L. Mahadevan, P. Koumoutsakos, Controlled gliding and perching through deep-
895 reinforcement-learning, *Phys. Rev. Fluids* 4 (2019) 093902.
- 896 [26] G. Novati, S. Verma, D. Alexeev, D. Rossinelli, W. M. van Rees, P. Koumoutsakos, Syn-
897 chronisation through learning for two self-propelled swimmers, *Bioinspir. Biomim.* 12 (2017)
898 036001.
- 899 [27] S. Verma, G. Novati, P. Koumoutsakos, Efficient collective swimming by harnessing vortices
900 through deep reinforcement learning, *Proc. Natl. Acad. Sci. U.S.A.* 115 (2018) 5849–5854.
- 901 [28] K. Lee, S. Kim, J. Choi, Deep reinforcement learning in continuous action spaces: a case study
902 in the game of simulated curling, in: *Procs. of the 35th International Conference on Machine*
903 *Learning*, 2018, pp. 4587–4596.
- 904 [29] X. Yan, J. Zhu, M. Kuang, X. Wang, Aerodynamic shape optimization using a novel optimizer
905 based on machine learning techniques, *Aerosp. Sci. Technol.* 86 (2019) 826–835.
- 906 [30] J. Viquerat, J. Rabault, A. Kuhnle, H. Ghraieb, E. Hachem, Direct shape optimization through
907 deep reinforcement learning, arXiv preprint [arXiv:1908.09885](https://arxiv.org/abs/1908.09885) (2019).
- 908 [31] P. Ma, Y. Tian, Z. Pan, B. Ren, D. Manocha, Fluid directed rigid body control using deep
909 reinforcement learning, *ACM Transactions on Graphics (TOG)* 37 (4) (2018) 1–11.
- 910 [32] L. Biferale, F. Bonaccorso, M. Buzicotti, P. Clark Di Leioni, K. Gustavsson, Zermelo’s prob-
911 lem: Optimal point-to-point navigation in 2D turbulent flows using reinforcement learning,
912 *Chaos* 29 (2019) 103138.
- 913 [33] J. Rabault, M. Kuchta, A. Jensen, U. Réglade, N. Cerardi, Artificial neural networks trained
914 through deep reinforcement learning discover control strategies for active flow control, *Journal*
915 *of Fluid Mechanics* 865 (2019) 281–302.
- 916 [34] F. Ren, H. Hu, H. Tang, Active flow control using machine learning: A brief review, *J.*
917 *Hydrodynam.* 32 (2020) 247–253.

- 918 [35] H. Tang, J. Rabault, A. Kuhnle, Y. Wang, T. Wang, Robust active flow control over a range
919 of Reynolds numbers using an artificial neural network trained through deep reinforcement
920 learning, *Phys. Fluids* 32 (2020) 053605.
- 921 [36] R. Paris, R. Beneddine, J. Dandois, Robust flow control and optimal sensor placement using
922 deep reinforcement learning, arXiv preprint arXiv:2006.11005 (2020).
- 923 [37] H. Xu, W. Zhang, J. Deng, J. Rabault, Active flow control with rotating cylinders by an
924 artificial neural network trained by deep reinforcement learning, *J. Hydrodynam.* 32 (2020)
925 254–258.
- 926 [38] D. Fan, L. Yang, Z. Wang, M. S. Triantafyllou, G. E. Karniadakis, Reinforcement learning for
927 bluff body active flow control in experiments and simulations, *Proc. Natl. Acad. Sci. U.S.A.*
928 117 (2020) 26091–26098.
- 929 [39] G. Beintema, A. Corbetta, L. Biferale, F. Toschi, Controlling Rayleigh-Bénard convection
930 via reinforcement learning, arXiv preprint arXiv:2003.14358 (2020).
- 931 [40] H. Kazmi, F. Mehmood, S. Lodeweyckx, J. Driesen, Gigawatt-hour scale savings on a budget
932 of zero: Deep reinforcement learning based optimal control of hot water systems, *Energy* 144
933 (2018) 159–168.
- 934 [41] T. Zhang, J. Luo, P. Chen, J. Liu, Flow rate control in smart district heating systems using
935 deep reinforcement learning, arXiv preprint arXiv:1912.05313 (2019).
- 936 [42] H. Ghraieb, P. Meliga, J. Viquerat, E. Hachem, Single-step deep reinforcement learning for
937 open-loop control of laminar and turbulent flows, *Phys. Rev. Fluids* (in revision) (2020).
- 938 [43] C. Gruau, T. Coupez, 3d tetrahedral, unstructured and anisotropic mesh generation with
939 adaptation to natural and multidomain metric, *Comput. Methods Appl. Mech. Engrg.* 194
940 (2005) 4951–4976.
- 941 [44] M. Bernacki, Y. Chastel, T. Coupez, R. Logé, Level set framework for the numerical modelling
942 of primary recrystallization in polycrystalline materials, *Scr. Mater.* 58 (2008) 1129–1132.
- 943 [45] Y. Mesri, H. Digonnet, T. Coupez, Advanced parallel computing in material forming with
944 CIMLib, *Eur. J. Comput. Mech.* 18 (2009) 669–694.
- 945 [46] T. Coupez, Metric construction by length distribution tensor and edge based error for
946 anisotropic adaptive meshing, *J. Comput. Phys.* 230 (2011) 2391–2405.
- 947 [47] S. Patankar, *Numerical heat transfer and fluid flow*, Taylor and Francis, 1980.
- 948 [48] S. V. Patankar, A numerical method for conduction in composite materials, flow in irregular
949 geometries and conjugate heat transfer, in: *Procs. of the 6th International Heat Transfer*
950 *Conference, 1978*, pp. 297–302.
- 951 [49] T. J. R. Hughes, G. R. Feijóo, L. Mazzei, J.-B. Quincy, The variational multiscale method -
952 a paradigm for computational mechanics, *Comput. Methods Appl. Mech. Engrg.* 166 (1998)
953 3–24.
- 954 [50] R. Codina, Stabilization of incompressibility and convection through orthogonal sub-scales in
955 finite element methods, *Comput. Methods Appl. Mech. Engrg.* 190 (2000) 1579–1599.
- 956 [51] Y. Bazilevs, V. M. Calo, J. A. Cottrell, T. J. R. Hughes, A. Reali, G. Scovazzi, Variational
957 multiscale residual-based turbulence modeling for large eddy simulation of incompressible
958 flows, *Comput. Methods Appl. Mech. Engrg.* 197 (2007) 173–201.
- 959 [52] E. Hachem, B. Rivaux, T. Kloczko, H. Digonnet, T. Coupez, Stabilized finite element method
960 for incompressible flows with high Reynolds number, *J. Comput. Phys.* 229 (23) (2010) 8643–
961 8665.

- 962 [53] E. Hachem, T. Kloczko, H. Dignonnet, T. Coupez, Stabilized finite element solution to handle
963 complex heat and fluid flows in industrial furnaces using the immersed volume method, *Int.*
964 *J. Numer. Methods Fluids* 68 (2012) 99–121.
- 965 [54] R. Codina, Stabilized finite element approximation of transient incompressible flows using
966 orthogonal subscales, *Comput. Methods Appl. Mech. Engrg.* 191 (2002) 4295–4321.
- 967 [55] E. Hachem, S. Feghali, R. Codina, T. Coupez, Immersed stress method for fluid-structure
968 interaction using anisotropic mesh adaptation, *Int. J. Numer. Meth. Eng.* 94 (2013) 805–825.
- 969 [56] R. Codina, Comparison of some finite element methods for solving the diffusion-convection-
970 reaction equation, *Comput. Methods Appl. Mech. Engrg.* 156 (1998) 185–210.
- 971 [57] S. Badia, R. Codina, Analysis of a stabilized finite element approximation of the transient
972 convection-diffusion equation using an ALE framework, *SIAM Journal on Numerical Analysis*
973 44 (2006) 2159–2197.
- 974 [58] A. N. Brooks, T. Hughes, Streamline upwind/Petrov-Galerkin formulations for convection
975 dominated flows with particular emphasis on the incompressible Navier-Stokes equations,
976 *Comput. Methods Appl. Mech. Engrg.* 32 (1982) 199–259.
- 977 [59] A. Galeão, E. Do Carmo, A consistent approximate upwind Petrov-Galerkin method for
978 convection-dominated problems, *Comput. Methods Appl. Mech. Engrg.* 68 (1988) 83–95.
- 979 [60] E. Hachem, H. Dignonnet, E. Massoni, T. Coupez, Immersed volume method for solving natural
980 convection, conduction and radiation of a hat-shaped disk inside a 3d enclosure, *International*
981 *Journal of numerical methods for heat & fluid flow* (2012).
- 982 [61] E. Hachem, H. Dignonnet, E. Massoni, T. Coupez, Immersed volume method for solving natural
983 convection, conduction and radiation of a hat-shaped disk inside a 3d enclosure, *Int. J. Numer.*
984 *Method. H.* 22 (2012) 718–741.
- 985 [62] E. Hachem, G. Jannoun, J. Veysset, M. Henri, R. Pierrot, I. Poitroult, E. Massoni, T. Coupez,
986 Modeling of heat transfer and turbulent flows inside industrial furnaces, *Simul. Model. Pract.*
987 *Th.* 30 (2013) 35–53.
- 988 [63] I. Goodfellow, Y. Bengio, A. Courville, *The Deep Learning Book*, MIT Press, 2017.
- 989 [64] R. S. Sutton, A. G. Barto, *Reinforcement learning: An introduction*, MIT press, 2018.
- 990 [65] A. Kakade, A natural policy gradient, *Adv. Neural Inf. Process Syst.* 14 (2001) 1531–1538.
- 991 [66] J. Schulman, S. Levine, P. Moritz, M. I. Jordan, P. Abbeel, Trust Region Policy Optimization,
992 *arXiv e-prints* (Feb. 2015). [arXiv:1502.05477](https://arxiv.org/abs/1502.05477).
- 993 [67] Y. Wang, H. He, X. Tan, Y. Gan, Trust region-guided proximal policy optimization, *arXiv*
994 *preprint arXiv:1901.10314* (2019).
- 995 [68] A. Hill, A. Raffin, M. Ernestus, A. Gleave, A. Kanervisto, R. Traore, P. Dhariwal, C. Hesse,
996 O. Klimov, A. Nichol, M. Plappert, A. Radford, J. Schulman, S. Sidor, Y. Wu, Stable baselines,
997 <https://github.com/hill-a/stable-baselines> (2018).
- 998 [69] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, W. Zaremba,
999 *Openai gym* (2016). [arXiv:arXiv:1606.01540](https://arxiv.org/abs/1606.01540).
- 1000 [70] G. de Vahl Davis, I. Jones, Natural convection in a square cavity: a comparison exercise, *Int.*
1001 *J. Numer. Methods Fluids* 3 (1983) 227–248.
- 1002 [71] H. Dixit, V. Babu, Simulation of high Rayleigh number natural convection in a square cavity
1003 using the lattice Boltzmann method, *Int. J. Heat Mass Transfer* 49 (2006) 727–739.
- 1004 [72] N. Markatos, K. Pericleous, Laminar and turbulent natural convection in an enclosed cavity,
1005 *Int. J. Heat Mass Transfer* 27 (1984) 772–775.

- 1006 [73] G. Barakos, E. Mitsoulis, Natural convection flow in a square cavity revisited: laminar and
1007 turbulent models with wall functions, *Int. J. Numer. Methods Fluids* 18 (1994) 695–719.
- 1008 [74] K. Khanafer, K. Vafai, M. Lightstone, Buoyancy-driven heat transfer enhancement in a two-
1009 dimensional enclosure utilizing nanofluids, *Int. J. Heat Mass Transfer* 46 (2003) 3639–3653.
- 1010 [75] A. Lazaric, M. Restelli, A. Bonarini, Reinforcement learning in continuous action spaces
1011 through sequential Monte Carlo methods, in: *Procs. of the 35th International Conference*
1012 *on Machine Learning*, 2018, pp. 4587–4596.
- 1013 [76] J. Sari, F. Cremonesi, M. Khalloufi, F. Cauneau, P. Meliga, Y. Mesri, E. Hachem, Anisotropic
1014 adaptive stabilized finite element solver for rans models, *Int. J. Numer. Methods Fluids* 86
1015 (2018) 717–736.
- 1016 [77] G. Guiza, A. Larcher, A. Goetz, L. Billon, P. Meliga, E. Hachem, Anisotropic boundary layer
1017 mesh generation for reliable 3D unsteady RANS simulations, *Finite Elem. Anal. Des.* 170
1018 (2020) 103345.
- 1019 [78] P. Meliga, E. Hachem, Time-accurate calculation and bifurcation analysis of the incompressible
1020 flow over a square cavity using variational multiscale modeling, *J. Comput. Phys.* 376 (2019)
1021 952–972.