

# Hierarchical generative modelling for autonomous robots

Kai Yuan (✉ [kai.yuan@ed.ac.uk](mailto:kai.yuan@ed.ac.uk))

The University of Edinburgh

Noor Sajid

University College London

Karl Friston

University College London <https://orcid.org/0000-0001-7984-8909>

Zhibin Li

The University of Edinburgh

---

## Article

**Keywords:** autonomous robotic operations, hierarchical generative modelling, autonomous task completion

**Posted Date:** November 15th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-965852/v1>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# Hierarchical generative modelling for autonomous robots

Kai Yuan<sup>1,3</sup>, Noor Sajid<sup>2,3</sup>, Karl Friston<sup>2</sup> and Zhibin Li<sup>1\*</sup>

<sup>1</sup>Edinburgh Centre for Robotics, The University of Edinburgh, UK

<sup>2</sup>Wellcome Centre for Human Neuroimaging, Queen Square Institute of Neurology, University College London, UK

<sup>3</sup>These authors contributed equally

\* Corresponding author: [zhibin.li@ed.ac.uk](mailto:zhibin.li@ed.ac.uk)

## ABSTRACT

Humans can produce complex movements when interacting with their surroundings. This relies on the planning of various movements and subsequent execution. In this paper, we investigated this fundamental aspect of motor control in the setting of autonomous robotic operations. We consider hierarchical generative modelling—for autonomous task completion—that mimics the deep temporal architecture of human motor control. Here, temporal depth refers to the nested time scales at which successive levels of a forward or generative model unfold: for example, the apprehension and delivery of an object requires both a global plan that contextualises the fast coordination of multiple local limb movements. This separation of temporal scales can also be motivated from a robotics and control perspective. Specifically, to ensure versatile sensorimotor control, it is necessary to hierarchically structure high-level planning and low-level motor control of individual limbs. We use numerical experiments to establish the efficacy of this formulation and demonstrate how a humanoid robot can autonomously solve a complex task requiring locomotion, manipulation, and grasping, using a hierarchical generative model. In particular, the humanoid robot can retrieve and deliver a box, open and walk through a door to reach the final destination. Our approach, and experiments, illustrate the effectiveness of using human-inspired motor control algorithms, which provide a scalable hierarchical architecture for autonomous performance of complex goal-directed tasks.

# 31 1. INTRODUCTION

---

32 Humans can control their bodies to produce intricate motor behaviours that are aligned with their  
33 objectives, e.g., opening the door or moving a box. These tasks rely on the coordination of two distinct  
34 processes: motor planning and execution <sup>1-4</sup>. To realise this coordination, human motor control unfolds  
35 at nested time scales at different levels of the neuronal hierarchy <sup>5,6</sup>, e.g., a high-level plan to arrive at a  
36 particular place can entail multiple, individual, reflexive low-level limb movements for walking.  
37 Similar distinctions have been introduced in the robotics and control literature <sup>7</sup>, with promising results  
38 of combining robot control with learnt motor movements. However, control and planning are often  
39 designed manually and separately, with little or no feedback between the two <sup>8</sup>. This structural  
40 separation limits the performance of the robot and places greater demands on human involvement at the  
41 deployment stage, even under controlled conditions.

42 Given these principles, we postulate that robotics systems should consider motor control as a  
43 consequence or outcome of a (learnt) hierarchical generative model, which can optimise and adapt  
44 future actions in uncertain environments. This proposal inherits from hierarchical formulations of motor  
45 control and ensuing planning or control as (active) inference <sup>9-21</sup>. Briefly, hierarchical generative models  
46 are a description of how sensory observations are generated, i.e., encodings of sensorimotor  
47 relationships relevant for motor control <sup>22,23</sup>. This implies that autonomous systems can be  
48 conceptualised in terms of hierarchical motor control, capable of dealing with complex and diverse  
49 tasks <sup>24</sup>.

50 Practically, hierarchical generative models speak to a spatiotemporal separation of planning and motor  
51 control that assures functional integration <sup>23,25,26</sup>. The requisite architectures can be regarded as a series  
52 of distinct levels that provide appropriate motor control <sup>27</sup>(Figure 1). The lowest level predicts the  
53 proprioceptive signals—generated using a forward model of the mechanics—and the kinetics that  
54 undergirds motor execution. This kinetics can be regarded as realising and equilibrium position or  
55 desired set point, without explicit modelling of task dynamics (c.f., the equilibrium point hypothesis)  
56 <sup>28,29</sup>. The level above generates the necessary sequence of fixed points, that are realised by the lower

57 level. This sequence speaks to the stability control that a human has over limbs, to perambulate in an  
58 upright manner over, say, a centre of gravity. The highest level then pertains to planning<sup>25,30</sup>. Here,  
59 different states represent endpoints of an agent's plan, e.g., move a box from a table to another.

60 To verify this proposition, we introduce a hierarchal generative model for autonomous robotic  
61 operations using a learning-based scheme. Our model has three distinct levels for planning and motor  
62 generation. The first rests on the equilibrium point hypothesis<sup>1</sup>; namely, a forward model that predicts  
63 the proprioceptive inputs reporting a desired movement, where reflexes realise the descending  
64 predictions or set points. Practically, the ensuing hierarchical model was optimised simultaneously at  
65 all levels. However, only the middle level planner had the access to state feedback, which allowed for  
66 a particular type of factorisation (i.e., functional specialisation) in our generative model. We leave  
67 further details for later sections.

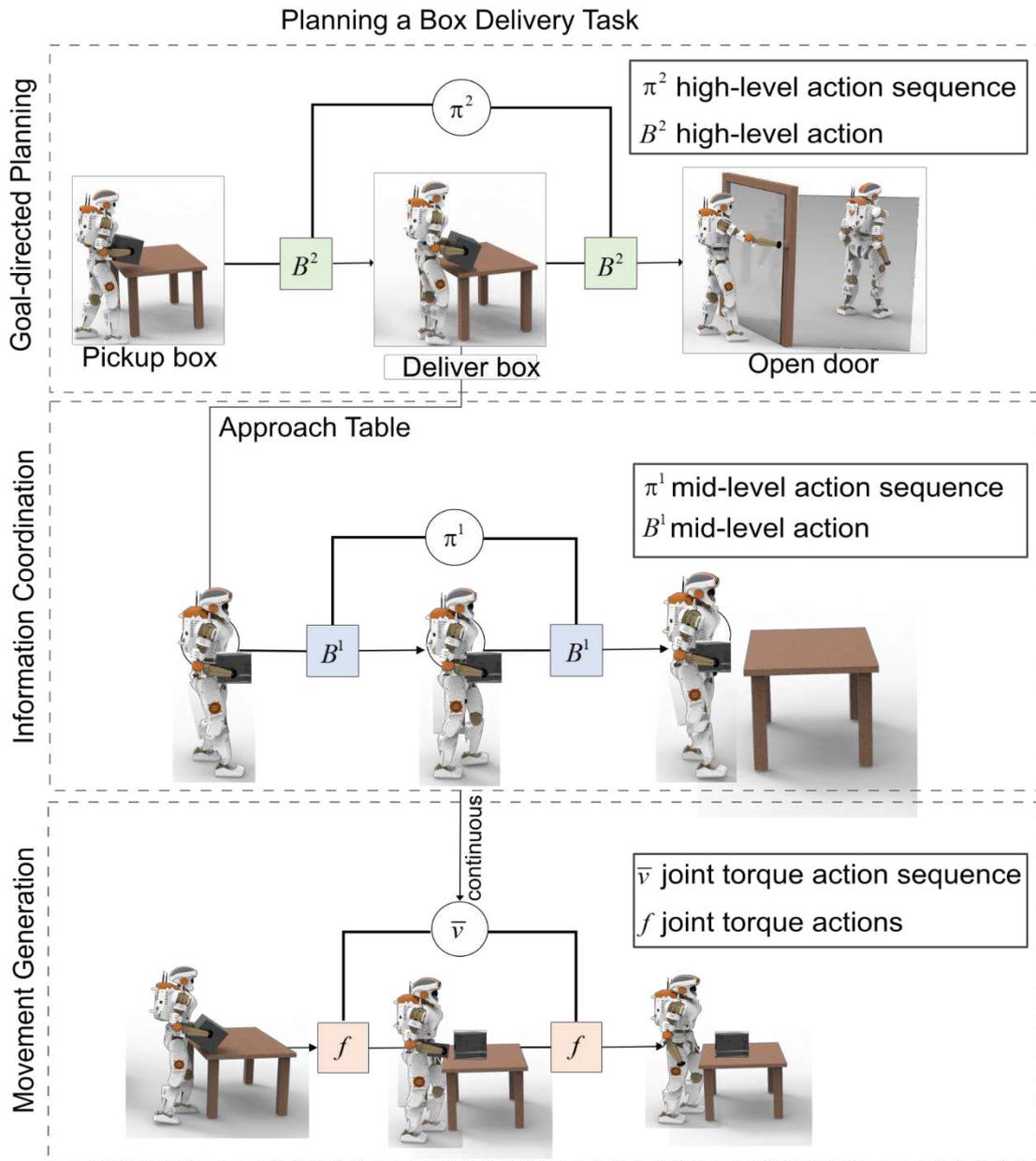
68 Using this generative model, we focused on solving tasks that comprised sequential, conditional  
69 decision-making in two scenarios. In the first scenario, the robot needed to reach a goal location by  
70 opening a closed door. The door only opened when a particular button was pressed. Here, the planner  
71 needed to coordinate both its legs and arms to (i) walk towards the door, (ii) and upon reaching it to  
72 move its arms to press the button, and (iii) finally walk through the door to reach the goal location.  
73 Conversely, the second scenario required the robot to carry a box from one table to another. Here, the  
74 planner needed to prescribe the sequence of approaching and picking up the box, walking towards the  
75 second table, and placing the box upon arrival. For both scenarios, our algorithm could perform coherent  
76 locomotion, manipulation, and grasping movements similar to humans, and therefore successfully solve  
77 these complex tasks.

78 The paper is organised as follows. In Section 2, we describe the realisations of two hierarchical  
79 generative models for motor control: human motor control and our robotic system capable of  
80 autonomous operations. The robotic system instantiates an implicit hierarchical generative model for  
81 motor control, which entails a bidirectional propagation of information between different levels of the

---

<sup>1</sup> Briefly, the equilibrium point hypothesis states that all movements are generated by the nervous system through a gradual transition of equilibrium points along a desired trajectory<sup>28,29</sup>.

82 generative model. Next, we show that the ensuing scheme can perform tasks remarkably similar to  
 83 humans (Section 3). In Section 4, we discuss the effectiveness of our hierarchical generative model,  
 84 how it may benefit potential applications, such as humanoids in disaster response scenarios, and provide  
 85 an outlook for future work. Lastly, in Section 5, we provide details of our implementation.



86

87 Figure 1: Pictorial representation of a hierarchical generative model for moving boxes – a form of motor  
 88 control. A generative model represents the conditional dependencies between states and how they cause  
 89 outcomes. For simplicity, we express this as filled squares that denote actions, and circles that represent  
 90 action sequences. The key aspect of this model is its hierarchical structure that represents sequences of  
 91 action over time. Here, actions at higher levels generate the initial actions for lower levels—that then  
 92 unfold to generate a sequence of actions: c.f., associative chaining. Crucially, lower levels cycle over a

93 sequence for each transition of the level above. It is this scheduling that endows the model with a deep  
 94 temporal structure. Particularly, planning (first row; highest level) to “deliver the box” generates the  
 95 actions for the information coordination level (second row; middle level) i.e., “movement towards the  
 96 table”. This in turn determines the initial actions for movement generation (third row; lowest level) of  
 97 arms to ‘place the box’ on the table.

## 98 2. HIERARCHICAL GENERATIVE MODELS FOR MOTOR 99 CONTROL

100 We hypothesise that motor control is a natural outcome of hierarchical generative modelling (Figure 1),  
 101 in particular, generative models that include the consequences of action. Hierarchical generative models  
 102 are descriptions of how sensory observations are generated – expressed as a joint distribution over the  
 103 unobservable causes of observable sensory input. These models can be factorised to represent particular  
 104 conditional dependencies that formalise sensorimotor contingencies<sup>22,23</sup>. Here, causes of proprioceptive  
 105 input can be predicted by fitting the model to observed sensory data<sup>31</sup>. This provides an explainable  
 106 and interpretable specification of decision-making within the robot, namely, planning and control as  
 107 inference. Interestingly, having a generative model gives for free the five core principles of hierarchical  
 108 motor control introduced in Merel et al., 2019 (see Table 1).

109 **Table 1.** Summary of the key principles of hierarchical motor control<sup>24</sup>, with exemplar realisations in  
 110 human motor control and our robotic system. We omit the principle of modular objectives here (sub-  
 111 systems trained to optimise specific objectives distinct from the global task objective) because a  
 112 factorised generative model architecture leads to distinct factor specific objectives at each level in the  
 113 hierarchy.

Principle	Description	Hierarchical generative models	Human motor control	Our robotics system for autonomous operations
Information factorisation	Different information is processed by distinct sub-systems.	Factorised distribution of appropriate latent states within the generative model.	Different sensory signals are routed to different parts in the hierarchy, e.g., what and where streams. These neuronal pathways can be characterised as factorised states responsible for sub-systems.	Only task-relevant sensory signals are used by individual levels, with irrelevant states hidden across levels. This speaks to an explicit factorisation of sensory signals and which parts of the system have access to them.
Partial autonomy	Lower hierarchical levels can semi-autonomously	The result of factorising state space into multiple levels can independently accomplish sub-goals at a	Semi-autonomous coordination of joint movement at lower levels (i.e., brainstem	Full autonomy and stability guaranteed at individual levels. Explicitly, we introduce stable mid-level

	produce outputs with minimum input from levels above.	(relatively) fast temporal scale.	and spinal cord). These operate at a faster temporal scale and do not require continuous input for higher levels.	and low-level motions for random higher-level inputs. This ensures that lower levels can independently perform fast movements.
Amortised control	Re-execute appropriate behaviours rapidly using learnt movements.	Learnt probability distributions that parameterise this generative model can be used for amortised control. That allows for habitual control based on previously learnt distributions.	The cerebellum is responsible for amortised control of deliberative and goal-directed behaviours, evoking fast habitual control for repeated actions.	The system learnt policies (i.e., action-state mappings) that provide habitual control for rapidly re-executing appropriate actions.
Multi-joint coordination	Degenerate coupling of different components operating as a whole for motor control.	Result of state factorisations that introduce flexible mapping across and within each level.	Different neuronal ensembles have distinct influences e.g., the red nucleus controls movements of the arms. Much like factorised states, these neuronal ensembles come together to produce intricate movements.	The system is equipped with multiple sub-structures (or policy mappings) that are responsible for specific actuator movement. Together these come across, and within levels, to produce particular motor movements.
Temporal abstraction	Abstraction of time across hierarchical levels.	A feature of hierarchical generative models, where higher levels evolve slower than and constrain, the level below.	Different levels evolve at different temporal and spatial scales, with the primary motor cortex responsible for planning (slow timescale) and spinal cord responsible for generation (fast timescale)	The three levels of the system evolve at different temporal scales – much like any hierarchical generative model. The high-level planning is at a slow timescale, mid-level stability control at medium timescale, and low-level joint control at a fast timescale.

114

115 The above considerations speak to the relevance of hierarchical generative models for motor control. In  
 116 the following, we describe the realisations in two domains: human and robotics (Figure 2).

117 **GENERATIVE MODELS OF HUMAN MOTOR CONTROL**

118 The inversion of forward models—that underwrites human motor control—generates continuous  
 119 proprioceptive predictions at the lowest level and propagates information to the highest levels that are  
 120 responsible for planning. Higher levels then pass messages down the hierarchy to generate particular  
 121 movements, via reflexive realisation of proprioceptive predictions (Figure 2). At the lowest

122 (sensorimotor) level, these proprioceptive predictions are usually multisensory. For simplicity, we limit  
123 ourselves to motor control by considering predictions of primary afferent signals from muscles.

124 The generative model's lowest (and fast-evolving) level includes the spinal cord and the brainstem. The  
125 local (factorised) circuits in these regions are responsible for organising movements, given descending  
126 predictions from different higher levels. Exemplar projections, from upper motor neurons in the cortex,  
127 control movement indirectly via pathways that project to the brainstem motor control centres, which, in  
128 turn, project to the local organising circuits in the brainstem and the spinal cord. Accordingly, these  
129 areas are responsible for evaluating the discrepancy between the proprioceptive inputs (primary  
130 afferents) and descending predictions of these signals, to drive muscle contraction via classical motor  
131 reflexes and accompanying musculoskeletal mechanics <sup>25,32</sup>. It is generally thought that the above  
132 systems coordinate multiple joint movements semi-autonomously over time <sup>33,342</sup>. At this level, different  
133 neuronal ensembles have distinct, partially autonomous influences, e.g., the red nucleus controls  
134 movements of the arms.

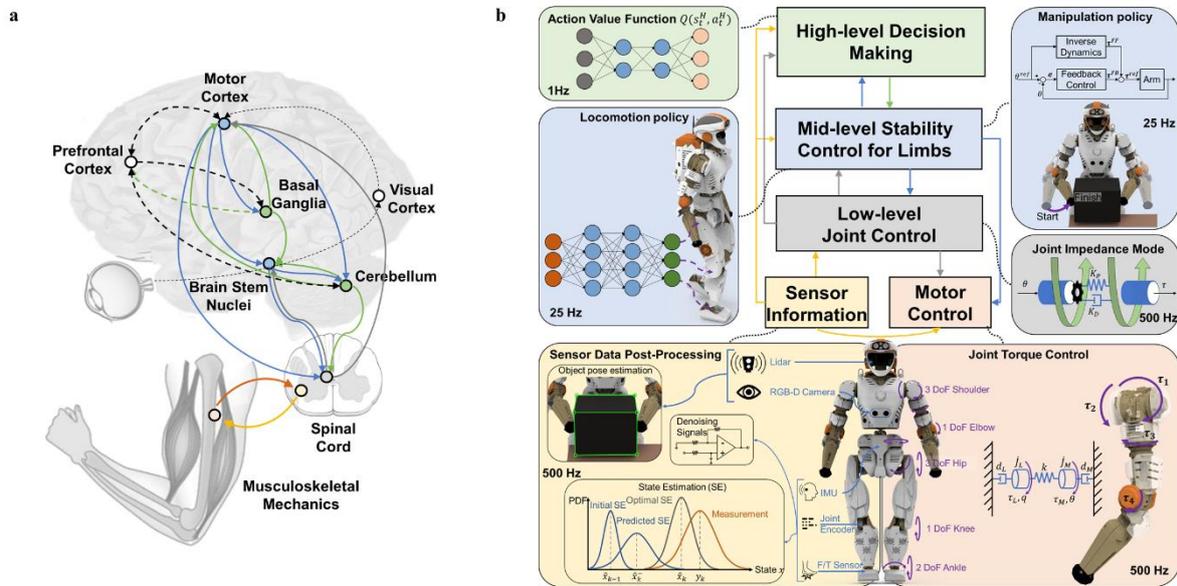
135 At intermediate levels of this generative model, one could consider the role of the cerebellum and the  
136 basal ganglia. The cerebellum receives ascending inputs from the spinal cord, and other areas, and  
137 integrates these to fine-tune motor activity. In other words, it does not initiate movement, but contributes  
138 to its coordination, precision, and speed, through fast non-deliberative mode of operation. Therefore, it  
139 can be thought of as being responsible for amortised (habitual) control of motor behaviour, which is  
140 characterised by subcortical and cortical interactions <sup>35-39</sup>. The cerebellum receives information from  
141 the motor cortex, processes this information, and sends motor impulses to skeletal muscles (via the  
142 spinal cord). Conversely, the basal ganglia is responsible for motor program selection, informed by the  
143 thalamus and the motor cortex <sup>40</sup>—appropriately orchestrating lower levels of the generative model  
144 (spinal cord via the brainstem). The basal ganglia are generally thought to support learning and action  
145 selection <sup>41</sup>.

---

<sup>2</sup> These trajectories are a succession of fixed points such that each fixed point is realised in a specified (e.g., theta) cycle over a particular time period. Then motor sequence can be regarded as filling in the gaps between a series of fixed points, where each fixed points is specified by the level above in very scheduled saltatory discrete fashion.

146 Higher levels of the generative model include the cerebral cortex and the thalamus. The thalamus  
 147 receives projections from the cerebellum (related to movement and sensory stimulation) and the basal  
 148 ganglia (related to wilful movements), and the thalamus projects directly to the primary motor and  
 149 premotor association cortices <sup>42</sup>. The cortex has access to factorised sensory streams of exteroceptive,  
 150 interoceptive and proprioceptive signals (e.g., visual, auditory, somatosensory, etc) and hence can  
 151 coordinate, contextualise or override habitual control elaborated in lower levels. Specifically, the  
 152 primary motor cortex is responsible for deliberative planning, control, and execution of voluntary  
 153 movements, for example, when learning a new motor skill prior to its habitisation or amortisation.

154



155

156 **Figure 2:** Algorithmic realisations of hierarchical control as inference. Panel A presents a schematic of  
 157 the generative model that underwrites human motor control, and Panel B depicts the implicit generative  
 158 model for a robotics system. The green nodes in Panel A and green boxes in Panel B refer to the highest  
 159 levels of human motor control and our implicit generative model respectively. In the generative model,  
 160 high-level decision making is realised as a neural network learned through Deep Reinforcement  
 161 Learning. The blue nodes in Panel A correspond to the middle level of human motor control and the  
 162 blue boxes in Panel B are intermediate level realisations, implemented as Deep Neural Network policy  
 163 learned through Deep Reinforcement Learning for Locomotion and Inverse Kinematics and Dynamics  
 164 policy for Manipulation. On the lowest level, depicted in grey nodes and boxes, a joint impedance  
 165 controller calculates the torques required for the actuation of the robot. Yellow and light red denote  
 166 sensor information and motor control respectively.

167

168 **IMPLICIT HIERARCHICAL GENERATIVE MODEL FOR A ROBOTICS SYSTEM**

169 Following the key principles of hierarchical motor control (Table 1) and the generative model in Figure  
170 1, we have constructed a generative model for a humanoid robot comprising three levels: high-level  
171 decision making, mid-level stability control, and low-level joint control (Figure 2). This hierarchical  
172 architecture rests on conditional independencies that result in factorised message passing between  
173 hierarchical levels<sup>3</sup>. For example, the mid-level locomotion policy only has access to the lower joint  
174 states and goals. Here, the temporal depth and structure motor planning is an outcome of specifying a  
175 hierarchical generative model, where level-specific policies are evaluated at different timescales. This  
176 involves each level assimilating (fast) evidence from the level below, in a way that is contextualised or  
177 selected by (slow) constraints, afforded by the level above.

178 In this example, the highest planning level, evolving at the slowest rate, selects an appropriate sequence  
179 of limb movements, which are needed to complete a particular sub-task. It decides where the hands  
180 should be and what direction to go. Practically, deep reinforcement learning (RL) is used to learn the  
181 correct generation of Cartesian position commands for the mid-level stability control.

182 These planning targets are realised at the level below that regulates the balance and stability of the robot  
183 during manipulation and locomotion. Manipulation is instantiated as a minimum-jerk model-predictive  
184 controller that moves the arms to the target positions provided by the high-level policy. Locomotion is  
185 implemented as a learnt policy, via deep RL, that coordinates legs—with twelve degrees of freedom—  
186 to reach the destination predicted by the higher level. Both of these kinds of policies are designed to  
187 ensure that infeasible set points—from the high level—are corrected for the mid-level stability control  
188 so that only stable joint target commands are supplied to the low-level joint controller.

189 Despite receiving inputs from other levels, each level has partial autonomy over its final predictions  
190 and goal. Furthermore, multi-joint coordination is realised by learning a policy that coordinates all joints

---

<sup>3</sup> It should be noted that the implementation used in our simulations is not formulated explicitly in terms of message passing or belief updating (which would usually be articulated in terms of Bayesian filtering for continuous states and belief propagation or variational message passing for discrete states). However, there exists an interpretation of the scheme in terms of expected states of the environment causing sensor inputs. Crucially, most of these causes correspond to the action of the robot itself.

191 of legs appropriately for the current state, while the arms coordinate their joints through Inverse  
192 Kinematics (IK).

193 The low-level joint controller is instantiated as joint impedance control and tracks the joint position  
194 commands afforded by the mid-level stability controller. Based on tuned stiffness and damping, the  
195 joint impedance control calculates the desired torque to attain target positions closely and smoothly.  
196 Lastly, the torque commands are tracked by the actuators, using embedded current control of onboard  
197 motor drivers.

### 198 **3. RESULTS**

---

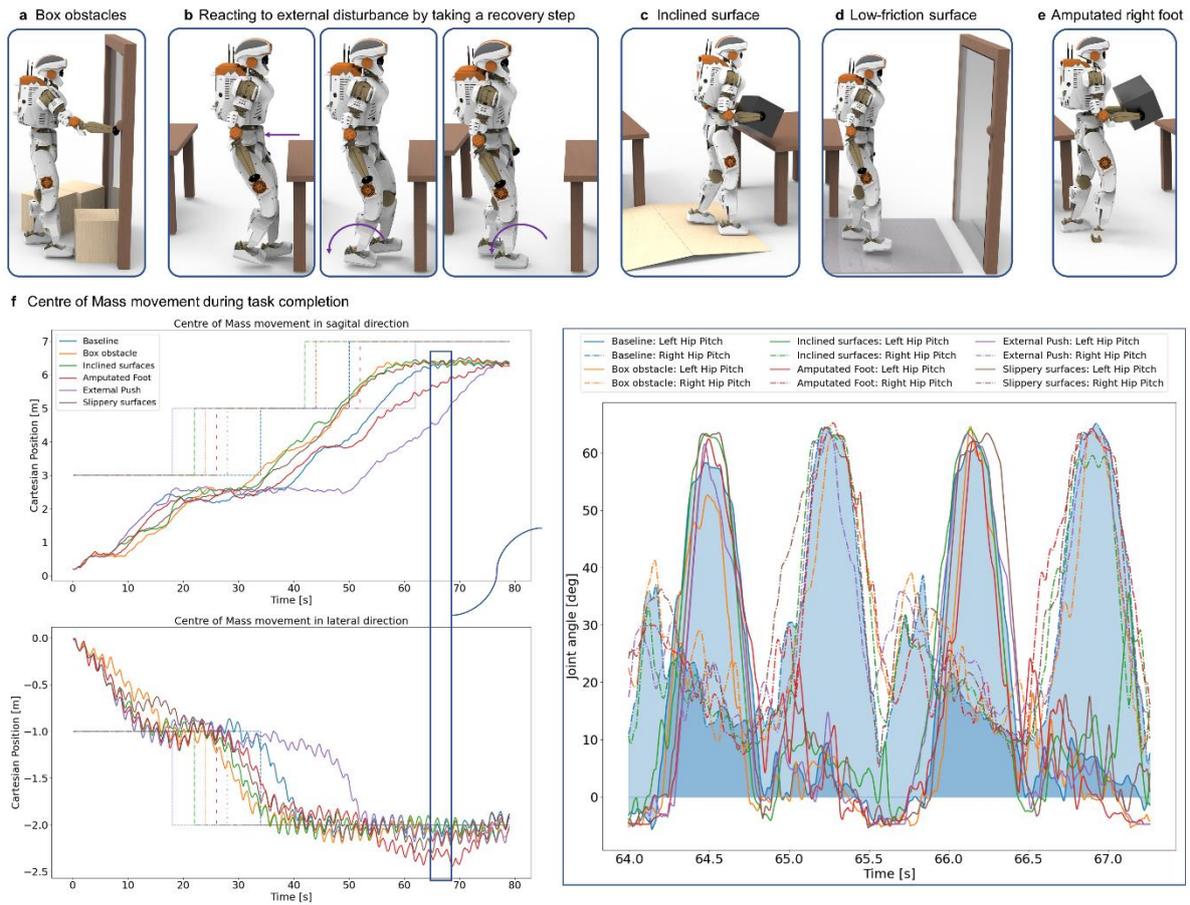
199 In the following, we illustrate how our proposed formulation of control architectures as inversion of a  
200 hierarchical generative model. Specifically, we show how this formulation enables a robot to learn how  
201 to complete a loco-manipulation task autonomously. Additionally, the learned policy features generality  
202 and robustness to uncertainty (Fig. 3 a-e), while evincing the core principles of hierarchical motor  
203 control. First, we validate the notion that hierarchical generative modelling is sufficient for autonomous  
204 robotic operations. We demonstrate this in two sequential decision-making tasks: moving a box from  
205 one table to another and opening a door by pressing a button.

206 Our implicit hierarchical generative model can successfully and autonomously achieve locomotion,  
207 manipulation, and grasping movements similar to humans, and solve all these complex tasks coherently  
208 and with internal consistency. Briefly, the highest policy level of the robot system determines the action  
209 sequence necessary for task completion and sends commands to the lower levels. This allows the robot  
210 to carry out the following actions: i) walk to the first table, ii) move arms to pick up the box, iii) walk  
211 to the second table, iv) and place the box on the table. Upon successful completion of this task, the same  
212 generative model can be used to perform the second task as well. This is achieved by using the high-  
213 level policy to open a closed door, instead of moving a box. The accompanying commands are sent to  
214 the lower levels responsible for limb stability and joint control. They instantiate mid-level locomotion

215 and manipulation policies to move to the door and position the body close enough to press and opening  
216 button. Having opened the door, the robot enters, and proceeds towards the final goal.

217 To assess the robustness and generality of this hierarchical scheme, we introduced several perturbations  
218 that were not encountered during training (Fig. 3). First, we introduced external perturbation by placing  
219 obstacles (i.e., 5kg box, Fig. 3a) in front of the robot and pushing its pelvis (Fig. 3b). The results were  
220 encouraging: the mid-level locomotion policy withstood both perturbations, moved the obstacle out of  
221 the way, and took a step to recover balance after the push. To test the performance further, we modified  
222 the environment with unseen conditions by adding a 5° inclined surface (Fig. 3c) and a low-friction  
223 glass plate (friction coefficient of 0.3, Fig. 3d) in front of the door. The robot could complete the task  
224 after each perturbation. More interestingly, we lesioned the robot by amputating its right foot (Fig. 3e).  
225 Despite this handicap (that was never encountered)—and with only a stump touching the ground in  
226 place of its right foot—our hierarchical control was sufficiently robust to deal with this situation and  
227 the robot was able to keep balance and complete the task.

228 Next, we evaluate whether the ensuing control architecture satisfies the key principles of hierarchical  
229 motor control (Table 1) that underwrite robust task performance.



230

231 **Figure 3:** Robustness of the system in light of perturbations. Panels A-E show how the robot completes  
 232 the task in perturbation test scenarios that it has not encountered during training and demonstrate the  
 233 robustness of our proposed method. From left to right, we place 5kg box-obstacles in front of the robot,  
 234 push the robot from the front, alter the floor with an inclined and slippery surface, and lesion the robot,  
 235 but removing the right foot. In Panel F, the amortised control is shown. The hip pitch joint motion is  
 236 used to show how the policy adapts to the perturbation and rapidly re-executes a motion to counteract  
 237 the perturbation.

### 238 INFORMATION FACTORISATION

239 In this system, factorisation exists across model levels and policy controls, each responsible for a  
 240 particular sort of information processing. This factorisation ensures that external perturbations have  
 241 minimum impact on task performance. For example, when disturbances involve robot legs, these  
 242 perturbations do not significantly impede task completion. This is because the perturbations are only  
 243 visible to the locomotion policy, and other levels and control policies are conditionally independent of  
 244 the disturbances and hence retain their nominal operations.

245 Practically, information factorisation necessitates clearly defined roles for each sub-system. Thus, any  
 246 failures in performance can be isolated and fine-tuned for future tasks. For example, if the robot falls

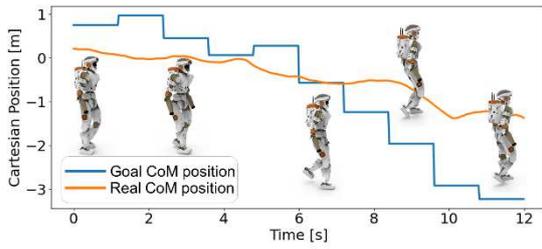
247 over while walking to a goal, the locomotion policy can be identified as the root cause, and hence  
248 improving the locomotion policy will resolve the issue without needing to modify the high-level planner  
249 or the manipulation policy. Further examples include oscillation of the robot limbs, which can be  
250 attributed to the low-level joint control; or walking in the wrong direction, which was due to the  
251 command from the high-level policy. From a theoretical perspective, factorisation of this sort  
252 corresponds to the structure of the generative model that can be decomposed into factors of a probability  
253 distribution. Almost universally, this results in certain conditional independencies that minimise the  
254 complexity of model inversion; namely, planning or control as inference<sup>9,11,13,43</sup>. This is important  
255 because it precludes over fitting and ensures generalisation. From a biological perspective, this kind of  
256 factorisation can be regarded as a functional segregation that is often associated with modular  
257 architectures and functional specialisation in the brain<sup>44</sup>.

## 258 **PARTIAL AUTONOMY**

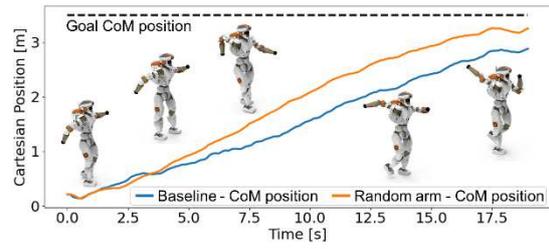
259 The system is designed with partial autonomy, i.e., minimum interference or support from other levels.  
260 Specifically, we implement a clear separation between the highest and intermediate levels, though they  
261 are learned together. This is particularly relevant because the high-level planning level could send  
262 unrealisable action sequences to the mid-level stability controller. Without partial autonomy, the robot  
263 might become unstable and unable to learn to move appropriately, given such random or potentially  
264 unstable high-level commands.

265 Figure 4 illustrates this case when the robot is provided with random commands to both the arms and  
266 legs. This causes the robot to walk in random directions (Fig. 4a) and the arms move around randomly  
267 (Fig. 4b). Despite imperfect motion tracking, the robot does not fall over and can complete the tasks  
268 despite incoherent intentions.

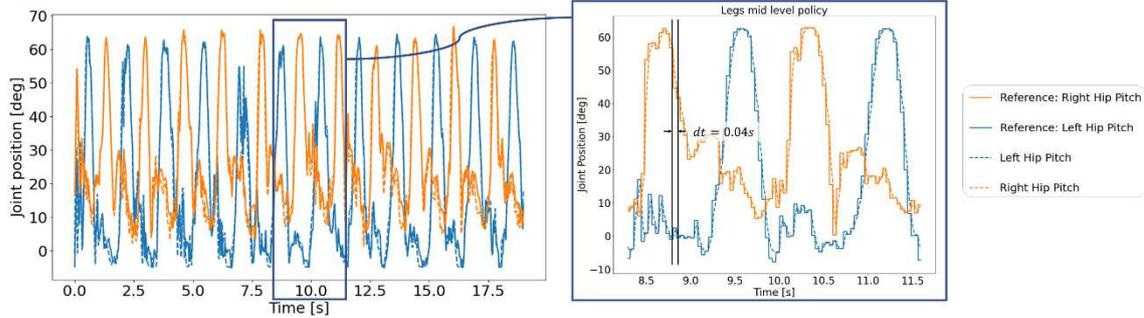
a Random leg motion: Sagittal CoM motion



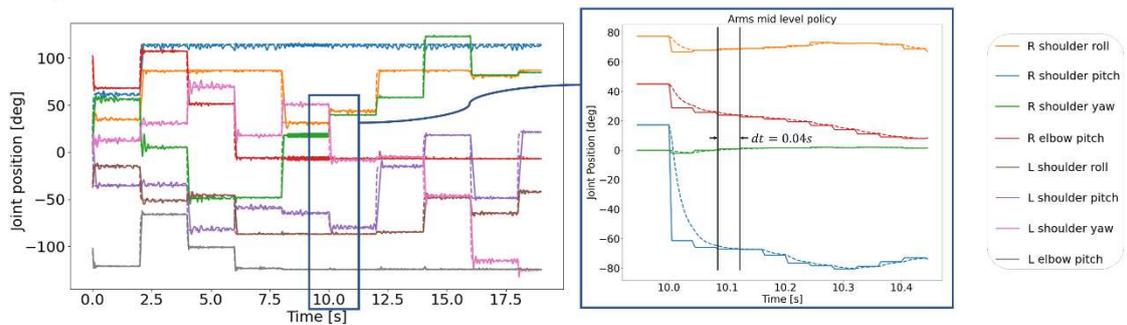
b Random arm motion: Sagittal CoM motion



c Leg joint movements



d Arm joint movements



269

270 **Figure 4:** State and temporal dynamics of the robot during task performance with random-high level  
 271 commands. Panels a and b show the sagittal motion of the Centre of Mass (CoM) while following  
 272 random leg and arm commands respectively. From the robot snapshots corresponding to the time  
 273 they're shown, the partial autonomy of the mid-level stability controllers can be seen, i.e., good  
 274 performance of the individual levels despite random and fast-changing command inputs. Panels c and  
 275 d show the leg and arm movements respectively. Here, the separation of temporal scales during planning  
 276 can be seen, where the high-level commands are provided at 0.5 Hz, and the mid-level commands are  
 277 realised at 25Hz. The joint commands are realised at 500Hz on the joint actuators. In the inset plots of  
 278 Panels c and d, it can be seen that the joint position trajectories evolve similarly as postulated in the  
 279 equilibrium point hypothesis.

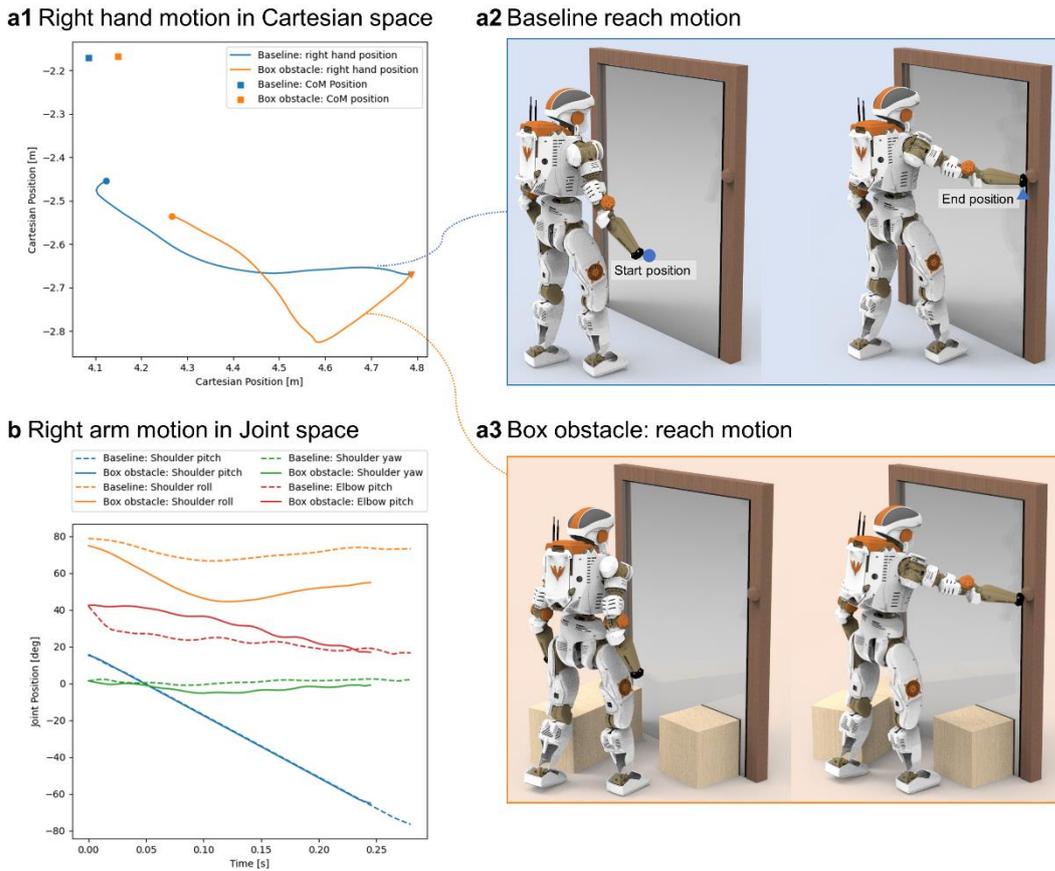
## 280 AMORTISED CONTROL

281 After training, the robot engages in amortised control with the ability to re-execute appropriate  
 282 behaviours rapidly using previously learnt movements. We observed this in the baseline and perturbed  
 283 task settings (inset trajectories in Fig. 3f), where the amortised locomotion policy was used to complete  
 284 the task without additional learning.

285 **MULTI-JOINT COORDINATION**

286 The robot has multiple sub-structures that are responsible for specific controls, and work together in  
287 different ways to generate particular motor movements. Figure 5a demonstrates this multi-joint  
288 coordination when pressing the button to open the door in the presence of an obstacle (Task 2). To  
289 achieve this, the right arm motions had to coordinate appropriately according to the initial hand position.  
290 Also, the shoulder roll (Fig. 5b orange line) and elbow (Fig. 5b red line) had to adjust and adapt  
291 differently from the baseline. Explicitly, these do not yield a fixed motion, instead, the manipulation  
292 policy coordinates these joints based on the Centre of Mass (CoM). Therefore, during the baseline  
293 reaching motion, the arms move differently than in the case of an obstructed box, where the CoM is in  
294 a different position (because boxes are obstructing the door).

295



296  
 297 **Figure 5:** Adaptive arm manipulations during opening the door task. Panel a shows the right arm’s  
 298 motion in cartesian space for the baseline (no obstacle) and the box obstacle scenario. When there is an  
 299 obstacle in front of the door, the arm motion must adapt, as the robot’s upper body position deviates  
 300 from its baseline position. The multi-joint coordination aspect can be seen in Panel b. The adapts to  
 301 the change of scenario, coordinating the multiple arm joints, such that the end position of the arm is pressing  
 302 the button to open the door, while being obstructed by box obstacles.

303

304 **TEMPORAL ABSTRACTION AND DEPTH**

305 By design, the three system levels evolve at different temporal scales. Figure 4 illustrates these distinct  
 306 scales as the robot perambulates. The highest policy level has a slow timescale of 0.5Hz (Fig. 4a). This  
 307 allows the lower levels to carry out the command in a partially autonomous way, i.e., uninterrupted.  
 308 Conversely, the mid-level stability control for limbs has a faster timescale at 25Hz (inset trajectories of  
 309 Fig. 4c and 4d). This is needed to generate rapid predictions for the locomotion and manipulation

310 policies. Finally, the low-level joint control executes these control commands at a frequency of 500Hz  
311 on the actuator level.

312

## 313 **4. DISCUSSION**

---

314 We have shown that hierarchical generative models can adequately equip agents with the ability to  
315 perform autonomously, in a context sensitive and robust fashion. Hierarchical generative modelling  
316 provides sensorimotor control systems with a transparent, and flexible, approach to interpret (and  
317 implement) robotic decision making. Importantly, this can be leveraged to improve task performance  
318 and identify system failures. It is worth acknowledging that the ensuing hierarchical framework adheres  
319 to the key principles of hierarchical motor control, and therefore offers the possibilities to (i) create a  
320 flexible system, (ii) rollout the policy and evaluate its performance, (iii) identify the root cause of the  
321 limitations in performance, and (iv) resolve issues in the root cause and improve the performance.

322 The proposed hierarchical generative model illustrates the effectiveness of hierarchical motor control  
323 principles for designing fully autonomous and robust robotics systems. We validated this approach on  
324 robot loco-manipulation tasks that required the retrieval and delivery of a box, opening a door, and  
325 walking towards a final goal. Furthermore, we showed the robustness of our system by demonstrating  
326 successful task completion despite previously unencountered perturbations, such as novel changes to  
327 the environment, unanticipated pushes, and even amputation of the right foot.

328 By providing robots with full autonomy with appropriate triage procedures, humans can be relieved of  
329 having to send low-level commands every few seconds (e.g., foot and hand placement) that require the  
330 interpretation of a steady stream of percept and sensor data—as in a shared autonomy and semi-  
331 autonomous paradigm. Consequently, human errors induced by a cognitive overload may be prevented  
332 in such robotics systems. Previously, these human errors have been observed in high-stress situations,  
333 despite being operated by experienced individuals with high social and hardware costs. A notable

334 example is the robot crash caused by incorrect foot placement during the DARPA Robotics Challenge—  
335 the most advanced competition to test the capabilities of anthropomorphic disaster-response robots.

336 Future work could evaluate the implications of more nuanced planning objectives, and their  
337 implications for such autonomous robotics systems. We expect that replacing our Q-learning planner,  
338 with more sophisticated objectives equipped to handle aleatoric and epistemic uncertainties (e.g.,  
339 expected free energy <sup>45-47</sup>), could improve the performance of the robot when dealing with volatile  
340 contingencies <sup>46</sup>. Furthermore, the robustness of our method could be evaluated through robotic  
341 neuropsychology <sup>48</sup>, i.e., introducing in-silico lesions to the robot system, and investigating their effect  
342 on the resulting policies, inference and behaviour.

343

## 344 **5. METHODS**

---

345 Here, we present the hardware implementation for inverting the (implicit) hierarchical generative (a.k.a.  
346 forward) model for autonomous robot control. First, we describe the robot platform on which the  
347 algorithms are implemented. Then, we detail the task that is solved autonomously by inverting the  
348 generative model; i.e., using the model to predict sensor inputs and using actuators to resolve the  
349 ensuing (proprioceptive) prediction errors. Next, we explain the generative model details including  
350 high-level decision making, mid-level stability control, and low-level joint control.

### 351 **ROBOT PLATFORM**

352 The motions for autonomous task completion are implemented on NASA’s humanoid Valkyrie<sup>49</sup>.  
353 Valkyrie was designed to operate in extra-terrestrial planetary space missions such as unmanned pre-  
354 deployment on Mars – it is 1.87m tall, consists of 44 Degrees of Freedom, and weighs 129kg with  
355 ranges of motion similar to humans. The 25 series-elastic actuators in arms, torso, and legs enable  
356 human-like locomotion and manipulation; all the joint limits are detailed in Figure S1. Valkyrie can  
357 sense the environment through proprioceptive and exteroceptive sensors, including a multitude of

358 gyroscopes, accelerometers, load cells, pressure sensors, sonar, LIDAR, depth cameras, and stereo  
359 sensors.

## 360 **TASKS OF INTEREST**

361 To demonstrate how in inversion of—or inference under—a hierarchical generative model solves  
362 complex tasks that require a particular sequence and coordination of locomotion and manipulation  
363 skills, we designed a task that demanded both coordination of limbs and reasoning about the sequence  
364 of actions. This task comprised four subtasks (Fig. S2): picking up a box from the first table, delivering  
365 the box to the second table, opening the door, and walking to the destination or goal position. To  
366 complete the task, all the subtasks had to be carried out in an exact sequence.

367 Our proposed framework allowed the robot to learn successful task completion through autonomous  
368 interactions with the environment. This was achieved by designing a reward (or utility) function for the  
369 high-level policy, such that cumulative maximisation of reward leads to task completion (see Section  
370 “High-level Decision Making”).

371 The trained policy learnt to first pick up the box, and hold the box using its hands, while approaching  
372 the second table. Once the box could be safely placed at the second table, the selected policy ensured  
373 the box was released and the robot moved towards the door, while keeping its arms in an appropriate  
374 position. Next, by pushing the button placed on the right side of the door, the robot could open the door.  
375 The learnt policy ensured that the final destination was approached at the right time, i.e., when the door  
376 was sufficiently open to walk through—towards the goal position.

## 377 **IMPLICIT GENERATIVE MODEL**

378 The generative model was realised by implementing three levels of control in a hierarchical manner  
379 (Fig. 2 in the manuscript): high-level decision making, mid-level stability control, and low-level joint  
380 control. All components were designed and trained separately, starting from the lowest level.

381 First, accurate and robust motor control needed to be guaranteed, such that the low-level joint position  
 382 control could be realised. Stiffness and damping parameters were tuned to track the references  
 383 accurately and compliantly, which provided the mid-level stability control. The mid-level stability  
 384 control consisted of a manipulation and a locomotion policy, which were individually designed. The  
 385 locomotion policy was trained to walk towards a commanded goal position, while the manipulation  
 386 policy was designed to place the hands on a target position. Finally, the high-level decision-making  
 387 policy was trained via Deep Reinforcement Learning, which learnt to provide appropriate commands  
 388 to these mid- and low-level policies.

### 389 **High-level Decision Making**

390 We achieved high-level decision making, the correct sequence and choices of robot actions, through  
 391 training a Deep Neural Network that approximated the action value function  $Q(s, a)$  over the  
 392 environment and choose the action  $a$  which yielded the highest value in state  $s$ .

### 393 *Double Q-learning*

394 We used Double Q-learning<sup>50</sup> to train a Q-network  $Q(s, a; \phi)$ , parametrised by weights  $\phi$ , to  
 395 approximate the true action value function  $Q(s, a)$ . At run-time, the action  $a$  was obtained as the  
 396 argument of the maximum Q-value  $a = \operatorname{argmax}_a Q(s, a; \phi)$  in state  $s$ . Compared to Deep Q-Learning<sup>51</sup>,  
 397 we used two separate Q-networks  $Q_1, Q_2$  for action selection and value estimation. Having two separate  
 398 Q-networks has previously shown to improve training stability<sup>50</sup>.

399 The network parameter  $\phi_i$  was obtained by  $\min_{\phi_i} L(\phi_i)$ :

$$400 \quad \min_{\phi_i} E \left[ \left( r + \gamma Q_j(s', a^*; \phi_j) - Q_i(s, a; \phi_i) \right)^2 \right],$$

401 with reward  $r$ , discount factor  $\gamma$ , network parameters  $\phi_i, \phi_j$ , Q-networks  $Q_i, Q_j$ , current state  $s$ , next  
 402 state  $s'$ , best action  $a^* = \operatorname{argmax}_a Q_i(s, a; \phi_i)$ . During training, either network parameters  $\phi_1$  or  $\phi_2$  was  
 403 randomly selected, trained, and used for action selection, while the other network parameter was used

404 to estimate the action value. The tuple  $(s, a, r, s') \sim U(D)$  was obtained from the Experience Replay  
 405 by uniformly sampling from buffer  $D$ , which was updated by online action rollout.

### 406 *Training Procedures*

407 The high-level policy sent and updated the actions  $a \in \mathcal{A} \subseteq \mathcal{R}^9$  at 1Hz frequency, which were the  
 408 positions in Cartesian space for the pelvis  $a_{\text{pelvis}} \in \mathcal{R}^3$ , left and right hands  $a_{\text{lh}}, a_{\text{rh}} \in \mathcal{R}^3$ . These actions  
 409  $a$  were executed by the mid-level stability controller. The locomotion policy constituted walking  
 410 towards the pelvis targets  $p_{\text{pelvis}}$ , while the manipulation policy placed the hands towards the hand  
 411 targets  $p_{\text{lhand}}, p_{\text{rhand}}$ .

412 The states  $s \in \mathcal{S} \subseteq \mathcal{R}^9$  were the vector  $\vec{s}_{\text{pelvis}} = a_{\text{pelvis}} - p_{\text{pelvis}} \in \mathcal{R}^3$  from current pelvis position  
 413  $p_{\text{pelvis}}$  to the pelvis target  $a_{\text{pelvis}}$ , the vector  $\vec{s}_{\text{hands}} = a_{\text{lh, rh}} - p_{\text{lh, rh}} \in \mathcal{R}^6$  from current hand positions  
 414  $p_{\text{lh, rh}}$  to left and right hand targets  $a_{\text{lh, rh}}$ . Lastly, three Boolean variables were provided as the  
 415 observation state when the door was open, the box was on the table, or the box was being carried.

416 The reward terms  $r_i$  were determined based on the task completion, such as whether the robot had  
 417 passed the delivery table, the arm joints were in the nominal position, the box was between the robot  
 418 hands, the box was at the delivery table, the door was open, and whether the robot was at the goal. The  
 419 weights  $w_i, i = 1, \dots, 6$  can be found in Figure S3.

420 At every time step, the reward  $r$  was the sum of sparse, Boolean states:

$$421 \quad r = w_1 r_{\text{pt}} + w_2 r_{\text{jn}} + w_3 r_{\text{bih}} + w_4 r_{\text{bot}} + w_5 r_{\text{do}} + w_6 r_{\text{ag}},$$

422 with passed table reward  $r_{\text{pt}}$ , joints nominal reward  $r_{\text{jn}}$ , box in hand reward  $r_{\text{bih}}$ , box on table reward  $r_{\text{bot}}$ ,  
 423 door open reward  $r_{\text{do}}$ , and at goal reward  $r_{\text{ag}}$ .

424 We terminated the episode early if the robot fell, or collided with itself, tables, or the door. Early  
 425 termination discouraged the policy from entering sub-optimal states (and actions) that would lead to  
 426 early termination as the cumulative reward becomes low, due to reaching the end of an episode.

427 We initialised the robot in different positions of the environment to allow the robot to encounter states  
428 that were hard to discover merely by exploration. Additionally, the robot was deliberately spawned at  
429 the configurations that were close to the final goal, in front of the open or closed door, and front of the  
430 tables.

### 431 **Mid-level Stability Control**

432 The mid-level stability control level consisted of two components: the manipulation policy for the arms  
433 realised as Model-Predictive Control (MPC) scheme, and a locomotion policy for the legs learned  
434 through Deep Reinforcement Learning.

#### 435 *Manipulation Policy*

436 The manipulation policy received Cartesian target positions for the hands from the High-level Decision-  
437 Making policy and sent joint position information to the Low-level Joint Controller. This was achieved  
438 by combining two parts (Fig. S4): Model-Predictive Control (MPC) that generated a stable, optimal  
439 trajectory in Cartesian space and Inverse Kinematics (IK)<sup>52</sup> that transformed desired actions from the  
440 Cartesian space to the joint space.

441 To provide the smoothest possible motions for the hands, we formulated the optimal control problem  
442 as the minimum-jerk optimisation, while satisfying dynamics constraints on the hands. The optimal  
443 trajectory was then implemented in an MPC fashion. The MPC control applied the first control input of  
444 the optimal input trajectory and then re-optimised based on the new state at the next control loop<sup>53</sup>. In  
445 this way, MPC successively solved an optimal control problem over a prediction horizon  $N$  and  
446 achieved feedback control, while ensuring optimality.

447 For the hand position  $x$ , an objective function  $J$  was designed to minimise jerk  $\ddot{x}$  (the input  $u$  of the  
448 system) with final time  $t_f$ :

$$449 \quad J = \frac{1}{2} \int_0^{t_f} \left( \frac{d^3 x(t)}{dt^3} \right)^2 dt = \frac{1}{2} \int_0^{t_f} u(t)^2 dt.$$

450 The Minimum Jerk MPC (MJMPC) solved the following constrained optimisation problem at every time  
 451 step at a frequency of 25Hz:

$$\begin{aligned}
 & 2 \min_{u(t)} && \frac{1}{2} \int_0^{t_f} u(t)^2 dt \\
 & \text{subject to} && \frac{d^3 x(t)}{dt^3} = u \\
 & && [x(0), \dot{x}(0), \ddot{x}(0)] = [x_0, \dot{x}_0, \ddot{x}_0] \\
 & && [x(t_f), \dot{x}(t_f), \ddot{x}(t_f)] = [x_f, \dot{x}_f, \ddot{x}_f] \\
 & && [x_{\min}, \dot{x}_{\min}, \ddot{x}_{\min}] \leq [x, \dot{x}, \ddot{x}] \leq [x_{\max}, \dot{x}_{\max}, \ddot{x}_{\max}],
 \end{aligned}$$

453 with initial condition  $[x_0, \dot{x}_0, \ddot{x}_0]$ , and terminal condition  $[x_f, \dot{x}_f, \ddot{x}_f]$ .

454 The resultant Cartesian trajectory  $p^d$  from MJMPC was transformed into joint position commands  $\theta^d$   
 455 through Inverse Kinematics (IK). More formally, IK described a transformation  $T: \mathcal{C} \rightarrow \mathcal{Q}$  from  
 456 Cartesian space  $\mathcal{C}$  to joint space  $\mathcal{Q}$ . The joint position commands  $\theta^d$  were then tracked by the low-level  
 457 joint position controller as described in Section “Low-level Joint Control”.

#### 458 ***Locomotion Policy***

459 The locomotion policy  $\pi(s; \theta)$  coordinated the 12 Degree of Freedom (DoF) leg joints and was  
 460 instantiated as a Deep Neural Network (network parameters  $\theta$ ) that received robot states  $s$  as inputs and  
 461 outputs 12 target joint positions for legs. It was trained through Soft-Actor Critic (SAC) <sup>54</sup>– an off-  
 462 policy Deep Reinforcement Learning (DRL) algorithm.

463 SAC optimised a maximum entropy objective  $J_{\text{SAC}(\pi)}$ :

$$464 \quad J_{\text{SAC}(\pi)} = \sum_{t=0}^T \mathbb{E} [r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot | s_t))],$$

465 with reward  $r$ , state  $s_t$  and action  $a_t$  at time  $t$ , temperature parameter  $\alpha$ , and policy entropy  $\mathcal{H}(\pi)$ . The  
 466 parameters  $\theta$  for policy  $\pi_\theta$  were obtained by minimising  $J_\pi(\theta)$ :

$$467 \quad J_\pi(\theta) = \mathbb{E} [\log \pi_\theta(a_t | s_t) - Q_\phi(s_t, a_t)].$$

468 The action-value function  $Q_\phi(s_t, a_t)$  was obtained by minimising the Bellman residual  $J_Q(\phi)$ :

469 
$$J_Q(\phi) = \mathbb{E} \left[ 1/2 \left( Q_\phi(s_t, a_t) - \hat{Q}(s_t, a_t) \right)^2 \right],$$

470 with Bellman equation  $\hat{Q}(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}[V_\psi(s_{t+1})]$  and discount factor  $\gamma$ . Estimation of the  
 471 value function  $V_\psi$  was obtained through minimising  $J_V(\psi)$ :

472 
$$J_V(\psi) = \mathbb{E} [ 1/2 \left( V_\psi(s_t) - \mathbb{E} [ Q_\theta(s_t, a_t) - \log \pi_\phi(a_t | s_t) ] \right)^2 ].$$

473 The training process of policy  $\pi_\theta$  is detailed in Algorithm 1.

---

**Algorithm 1** Pseudocode for SAC

---

1:  $\mathcal{D} \leftarrow \emptyset$ , initialise replay buffer  
 2: **for** iter=1,2,... **do**  
 3:     **while** collected samples < batch size **do**  
 4:         Sample and perform action  $a_t \sim \pi_\theta$   
 5:         Collect  $d = \{s_t, a_t, r(s_t, a_t), s_{t+1}\}$   
 6:         Store sample  $d$  in replay buffer  $\mathcal{D} \leftarrow \mathcal{D} \cup d$   
 7:     **for** policy update=1,...,K **do**  
 8:         Sample batch from replay buffer  $\mathcal{D}$  for update  
 9:         Update policy  $\pi_\theta$  via (1.5)  
 10:         Update action value function  $Q_\phi$  via (1.6)  
 11:         Update value function  $V_\psi$  via (1.7)

---

474

475 The training procedures, including the design of reward, action space, and state-space, are as in<sup>55</sup>. The  
 476 action space  $\mathcal{A} \in \mathcal{R}^{12}$  were the joint positions of the 12 Degree of Freedoms (DoF) of the legs (for each  
 477 leg 3 DoF hip, 1 DoF knee, 2 DoF ankle). The target joint positions were tracked by the low-level joint  
 478 controller (see Section “Low-level Joint Control”).

479 The state-space  $\mathcal{S} \in \mathcal{R}^{27}$  consisted of the desired pelvis position/reference (the walking destination),  
 480 and proprioceptive feedback states of the robot including pelvis orientation, linear and angular velocity  
 481 of the pelvis, force of both feet, joint positions of the legs, and the gait phase.

482 The reward comprised of an imitation term and a task term similar to :

483 
$$r_i = w_i r_{\text{imitation}} + w_t r_{\text{task}},$$

484 with weights  $w_i, w_t$  and reward terms  $r_{\text{imitation}}, r_{\text{task}}$  for imitation and task respectively.

485 The aim of  $r_{\text{imitation}}$  was to imitate the joint position, feet pose, and contact pattern of a reference motion  
486 capture trajectory as closely as possible. The reward term  $r_{\text{task}}$  rewarded upright posture, short distances  
487 to the goal position, and regularised the joint velocity and torque.

#### 488 **Low-level Joint Control**

489 The low-level joint control tracked the target joint positions  $\theta^d$  provided by the mid-level stability  
490 controller. It was realised as *Joint Impedance Controller* that regulated around the set point.  
491 Additionally, a feed-forward controller (*Inverse Dynamics* block in Fig. S5) provided joint torques  $\tau^{FF}$   
492 that compensated for the dynamical influences from gravity and Coriolis forces and hence achieved  
493 more accurate tracking of the desired joint motions  $\theta^d$ .

#### 494 **Feedback Control**

495 The desired control effort was calculated using stiffness  $K_{P_i}$  and damping gains  $K_{D_i}$ . The joint  
496 impedance control calculated the desired joint torque  $\tau^d$  using position  $\theta$  and its derivative  $\dot{\theta}$ :

$$497 \quad \tau^d = K_{P_1}(\theta^d - \theta) - K_{D_1}\dot{\theta}.$$

498 At the actuator level, the motor driver implemented an internal current control to track the desired joint  
499 torque  $\tau^d$  using a proportional-derivative law, where the desired motor current  $I$  was computed as:

$$500 \quad I = K_{P_2}(\tau^d - \tau) - K_{D_2}\dot{\tau}.$$

#### 501 **Feedforward Control**

502 The feedforward torques  $\tau^{FF}$  were computed through Inverse Dynamics, which transformed  $T: \theta \rightarrow T$   
503 from desired joint motions  $\theta^d \in \theta$  to joint torques  $\tau^{FF} \in T$ , based on the equations of motions  
504 satisfying the Newton–Euler dynamics of multi-link rigid bodies given the joint configuration  $\theta^d, \dot{\theta}^d$ .  
505 These resultant torques  $\tau^{FF}$  compensated for gravity and Coriolis forces given the desired joint  
506 configuration  $\theta^d$ .

507

- 
- 509 1 Li, N., Chen, T.-W., Guo, Z. V., Gerfen, C. R. & Svoboda, K. A motor cortex circuit for motor  
510 planning and movement. *Nature* **519**, 51-56, doi:10.1038/nature14178 (2015).
- 511 2 Rosenbaum, D. A., Vaughan, J., Barnes, H. J. & Jorgensen, M. J. Time course of movement  
512 planning: Selection of handgrips for object manipulation. *Journal of Experimental Psychology:*  
513 *Learning, Memory, and Cognition* **18**, 1058-1073, doi:10.1037/0278-7393.18.5.1058 (1992).
- 514 3 Svoboda, K. & Li, N. Neural mechanisms of movement planning: motor cortex and beyond.  
515 *Current opinion in neurobiology* **49**, 33-41 (2018).
- 516 4 Hogan, N. Planning and execution of multijoint movements. *Canadian Journal of Physiology*  
517 *and Pharmacology* **66**, 508-517 (1988).
- 518 5 Honey, C. J. *et al.* Slow cortical dynamics and the accumulation of information over long  
519 timescales. *Neuron* **76**, 423-434, doi:10.1016/j.neuron.2012.08.011 (2012).
- 520 6 Murray, J. D. *et al.* A hierarchy of intrinsic timescales across primate cortex. *Nat Neurosci* **17**,  
521 1661-1663, doi:10.1038/nn.3862 (2014).
- 522 7 Diedrichsen, J., Shadmehr, R. & Ivry, R. B. The coordination of movement: optimal feedback  
523 control and beyond. *Trends in cognitive sciences* **14**, 31-39, doi:10.1016/j.tics.2009.11.004  
524 (2010).
- 525 8 Johnson, M. *et al.* Team IHMC's Lessons Learned from the DARPA Robotics Challenge Trials.  
526 *Journal of Field Robotics* **32**, 192-208 (2015).
- 527 9 Attias, H. in *Proc. of the 9th Int. Workshop on Artificial Intelligence and Statistics*.
- 528 10 Baker, C. L., Saxe, R. & Tenenbaum, J. B. Action understanding as inverse planning. *Cognition*  
529 **113**, 329-349, doi:10.1016/j.cognition.2009.07.005 (2009).
- 530 11 Botvinick, M. & Toussaint, M. Planning as inference. *Trends Cogn Sci.* **16**, 485-488 (2012).
- 531 12 Maisto, D., Donnarumma, F. & Pezzulo, G. Divide et impera: subgoaling reduces the  
532 complexity of probabilistic inference and problem solving. **12**, 20141335,  
533 doi:10.1098/rsif.2014.1335 (2015).
- 534 13 Kaplan, R. & Friston, K. J. Planning and navigation as active inference. *Biological cybernetics*  
535 **112**, 323-343, doi:10.1007/s00422-018-0753-2 (2018).
- 536 14 Tani, J. & Nolfi, S. Learning to perceive the world as articulated: an approach for hierarchical  
537 learning in sensory-motor systems. *Neural Netw* **12**, 1131-1141 (1999).
- 538 15 Tani, J. Learning to generate articulated behavior through the bottom-up and the top-down  
539 interaction processes. *Neural Netw.* **16**, 11-23 (2003).
- 540 16 Jung, M., Hwang, J. & Tani, J. Self-Organization of Spatio-Temporal Hierarchy via Learning  
541 of Dynamic Visual Image Patterns on Action Sequences. *PLOS ONE* **10**, e0131214,  
542 doi:10.1371/journal.pone.0131214 (2015).
- 543 17 Matsumoto, T. & Tani, J. Goal-Directed Planning for Habituated Agents by Active Inference  
544 Using a Variational Recurrent Neural Network. *Entropy* **22**, 564, doi:10.3390/e22050564  
545 (2020).
- 546 18 Haruno, M., Wolpert, D. M. & Kawato, M. Hierarchical MOSAIC for movement generation.  
547 *International Congress Series* **1250**, 575-590, doi:[https://doi.org/10.1016/S0531-](https://doi.org/10.1016/S0531-5131(03)00190-0)  
548 [5131\(03\)00190-0](https://doi.org/10.1016/S0531-5131(03)00190-0) (2003).
- 549 19 Wolpert, D. M., Doya, K. & Kawato, M. A unifying computational framework for motor  
550 control and social interaction. *Philos Trans R Soc Lond B Biol Sci.* **358**, 593-602 (2003).
- 551 20 Morimoto, J. & Doya, K. Acquisition of stand-up behavior by a real robot using hierarchical  
552 reinforcement learning. *Robotics and Autonomous Systems* **36**, 37-51 (2001).
- 553 21 Doya, K. *Bayesian brain: Probabilistic approaches to neural coding.* (MIT press, 2007).
- 554 22 Baltieri, M. & Buckley, C. L. Generative models as parsimonious descriptions of sensorimotor  
555 loops. *Behavioral and Brain Sciences* **42**, e218, doi:10.1017/S0140525X19001353 (2019).
- 556 23 Friston, K. J., Parr, T. & de Vries, B. The graphical brain: Belief propagation and active  
557 inference. *Network neuroscience* **1**, 381-414, doi:10.1162/NETN\_a\_00018 (2017).
- 558 24 Merel, J., Botvinick, M. & Wayne, G. Hierarchical motor control in mammals and machines.  
559 *Nature Communications* **10**, 5489, doi:10.1038/s41467-019-13239-6 (2019).

560 25 Parr, T., Limanowski, J., Rawji, V. & Friston, K. The computational neurology of movement  
561 under active inference. *Brain : a journal of neurology* (2021).

562 26 Jackson, J. H. On certain relations of the cerebrum and cerebellum (on rigidity of hemiplegia  
563 and on paralysis agitans). *Brain : a journal of neurology* **22**, 621-630 (1899).

564 27 Pezzulo, G., Rigoli, F. & Friston, K. J. Hierarchical Active Inference: A Theory of Motivated  
565 Control. *Trends in cognitive sciences* **22**, 294-306, doi:10.1016/j.tics.2018.01.009 (2018).

566 28 Feldman, A. G. & Levin, M. F. The equilibrium-point hypothesis--past, present and future.  
567 *Advances in experimental medicine and biology* **629**, 699-726, doi:10.1007/978-0-387-77064-  
568 2\_38 (2009).

569 29 Perrier, P., Ostry, D. J. & Laboissière, R. The Equilibrium Point Hypothesis and Its Application  
570 to Speech Motor Control. *Journal of Speech, Language, and Hearing Research* **39**, 365-378,  
571 doi:doi:10.1044/jshr.3902.365 (1996).

572 30 Botvinick, M. M., Niv, Y. & Barto, A. C. Hierarchically organized behavior and its neural  
573 foundations: A reinforcement learning perspective. *Cognition* **113**, 262-280,  
574 doi:10.1016/j.cognition.2008.08.011 (2009).

575 31 Aitchison, L. & Lengyel, M. With or without you: predictive coding and Bayesian inference in  
576 the brain. *Current opinion in neurobiology* **46**, 219-227 (2017).

577 32 Adams, R. A., Shipp, S. & Friston, K. J. Predictions not commands: active inference in the  
578 motor system. *Brain Struct Funct.* **218**, 611-643 (2013).

579 33 Bizzi, E., Mussa-Ivaldi, F. & Giszter, S. Computations underlying the execution of movement:  
580 a biological perspective. *Science* **253**, 287-291, doi:10.1126/science.1857964 (1991).

581 34 Marder, E. & Bucher, D. Central pattern generators and the control of rhythmic movements.  
582 *Current biology : CB* **11**, R986-996, doi:10.1016/s0960-9822(01)00581-4 (2001).

583 35 Miall, R. C., Weir, D. J., Wolpert, D. M. & Stein, J. F. Is the cerebellum a smith predictor? *J*  
584 *Mot Behav.* **25**, 203-216 (1993).

585 36 Doya, K. What are the computations of the cerebellum, the basal ganglia and the cerebral cortex?  
586 *Neural Netw* **12**, 961-974, doi:10.1016/s0893-6080(99)00046-5 (1999).

587 37 Freeman, J. H. & Steinmetz, A. B. Neural circuitry and plasticity mechanisms underlying delay  
588 eyeblink conditioning. *Learning & memory (Cold Spring Harbor, N.Y.)* **18**, 666-677,  
589 doi:10.1101/lm.2023011 (2011).

590 38 Koziol, L. F. *et al.* Consensus paper: the cerebellum's role in movement and cognition.  
591 *Cerebellum (London, England)* **13**, 151-177, doi:10.1007/s12311-013-0511-x (2014).

592 39 Ramnani, N. Automatic and controlled processing in the corticocerebellar system. *Prog Brain*  
593 *Res* **210**, 255-285, doi:10.1016/b978-0-444-63356-9.00010-8 (2014).

594 40 Grillner, S., Wallén, P., Saitoh, K., Kozlov, A. & Robertson, B. Neural bases of goal-directed  
595 locomotion in vertebrates—an overview. *Brain research reviews* **57**, 2-12 (2008).

596 41 Jueptner, M., Frith, C., Brooks, D., Frackowiak, R. & Passingham, R. Anatomy of motor  
597 learning. II. Subcortical structures and learning by trial and error. *Journal of neurophysiology*  
598 **77**, 1325-1337 (1997).

599 42 Sommer, M. A. The role of the thalamus in motor control. *Current Opinion in Neurobiology*  
600 **13**, 663-670, doi:<https://doi.org/10.1016/j.conb.2003.10.014> (2003).

601 43 Millidge, B. Deep Active Inference as Variational Policy Gradients. *arXiv e-prints*,  
602 arXiv:1907.03876 (2019).

603 44 Parr, T., Sajid, N. & Friston, K. J. Modules or Mean-Fields? *Entropy* **22**, 552 (2020).

604 45 Friston, K., FitzGerald, T., Rigoli, F., Schwartenbeck, P. & Pezzulo, G. Active Inference: A  
605 Process Theory. *Neural Comput* **29**, 1-49, doi:10.1162/NECO\_a\_00912 (2017).

606 46 Sajid, N., Ball, P. J., Parr, T. & Friston, K. J. Active inference: demystified and compared.  
607 *Neural computation* **33**, 674-712 (2021).

608 47 Da Costa, L. *et al.* Active inference on discrete state-spaces: A synthesis. *Journal of*  
609 *Mathematical Psychology* **99**, 102447, doi:<https://doi.org/10.1016/j.jmp.2020.102447> (2020).

610 48 Sajid, N. *et al.* Simulating lesion-dependent functional recovery mechanisms. *Scientific Reports*  
611 **11**, 7475, doi:10.1038/s41598-021-87005-4 (2021).

612 49 Radford, N. A. *et al.* Valkyrie: Nasa's first bipedal humanoid robot. *Journal of Field Robotics*  
613 **32**, 397-419 (2015).

614 50 Hasselt, H. Double Q-learning. *Advances in neural information processing systems* **23**, 2613–  
615 2621 (2010).  
616 51 Mnih, V. *et al.* Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*  
617 (2013).  
618 52 Siciliano, B., Khatib, O. & Kröger, T. *Springer handbook of robotics*. Vol. 200 (Springer, 2008).  
619 53 Rawlings, J. B., Mayne, D. Q. & Diehl, M. *Model predictive control: theory, computation, and*  
620 *design*. Vol. 2 (Nob Hill Publishing Madison, WI, 2017).  
621 54 Haarnoja, T., Zhou, A., Abbeel, P. & Levine, S. in *International conference on machine*  
622 *learning*. 1861–1870.  
623 55 Yang, C., Yuan, K., Heng, S., Komura, T. & Li, Z. Learning natural locomotion behaviors for  
624 humanoid robots using human bias. *IEEE Robotics and Automation Letters* **5**, 2610–2617  
625 (2020).

## 626 AUTHOR CONTRIBUTIONS

---

627 K.Y. and Z.L. conceptualised the robot control architecture. K.Y., N.S., and K.F. designed the  
628 hierarchical generative model. K.Y. implemented the model and performed the robotic experiments.  
629 K.Y. and N.S. wrote the manuscript. All authors contributed to and edited the manuscript.

## 630 COMPETING INTERESTS

---

631 The authors declare no competing interests.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [supplementary.pdf](#)