

A Hybrid Controller based on the Egocentric Perceptual Principle

Zinovi Rabinovich*, Nicholas R. Jennings

Electronics and Computer Science, University of Southampton, Southampton, SO17 1BJ, UK

Abstract

In this paper we extend the control methodology based on Extended Markov Tracking (EMT) by providing the control algorithm with capabilities to calibrate and even partially reconstruct the environment's model. This enables us to resolve the problem of performance deterioration due to model incoherence, a problem faced in all model-based control methods. The new algorithm, *Ensemble Actions EMT (EA-EMT)*, utilises the initial environment model as a library of state transition functions and applies a variation of prediction with experts to assemble and calibrate a revised model. By so doing, this is the first hybrid control algorithm that enables on-line adaptation within the egocentric control framework which dictates the control of an agent's perceptions, rather than an agent's environment state. In our experiments, we performed a range of tests with increasing model incoherence induced by three types of exogenous environment perturbations: *catastrophic* – the environment becomes completely inconsistent with the model, *deviating* – some aspect of the environment behaviour diverges compared to that specified in the model, and *periodic* – the environment alternates between several possible divergences. The results show that EA-EMT resolved model incoherence and significantly outperformed its EMT predecessor by up to 95%.

Keywords: hybrid control, perceptual control, dynamics based control, Kullback-Leibler divergence

1. Introduction

2 Egocentric perceptual control (EPC) formulates a control problem in terms of an
3 agent's perceptions, i.e. its internal interpretation of sensory input, rather than the
4 actual environment state [1]. As a direct outcome of this representation, any task

*Corresponding author

Email addresses: zr@ecs.soton.ac.uk (Zinovi Rabinovich),
nrj@ecs.soton.ac.uk (Nicholas R. Jennings)

5 that an agent performs is expressed as a preference over perceptions, and the op-
6 timality criteria follows suit. In fact, from this egocentric point of view, changes
7 in the environment are simply a means to alter and control the agent's perceptions.
8 As a technical example consider instrument flight rules (IFR), the regulations and
9 procedures for flying aircraft by referring only to the aircraft instrument panel for
10 navigation. These rules describe the instrument readings that a pilot (and hence the
11 auto-pilot control algorithm) has to maintain, therefore referring to the perceived
12 flight parameters, rather than the factual physical state of the plane. Notice that the
13 instrument readings are indeed *perceptions*, the interpretations of the automated
14 sensors, rather than the observations or measurements that they make. To date,
15 EPC has been used in a variety of domains, including sensory-based navigation
16 of autonomous robots, where all the necessary information is represented through
17 perceptions, such as maps or landmarks (see e.g. [2, 3]). In fact, one of the most
18 successful control approaches in robotics, the behaviour-based control (BBR) [4],
19 can be seen to be a particular instantiation of the EPC. In more detail, in BBR a
20 complex behaviour with desired properties is obtained by means of arbitration and
21 fusion of a set of simple mappings (basic behaviours) from perceptions to actions¹.
22 Starting from the simplest basic behaviours, that are enacted once some key per-
23 ception is formed, and ending with complex arbitration of a BBR scheme, all key
24 features of decision making are based on perceptual information, therefore con-
25 forming BBR to the EPC view. Moreover, EPC is inherent to behaviour patterns
26 found in nature or based on human intuition and psychology (e.g. [5, 6] and refer-
27 ences therein). It enables, for instance, a quick design of individual behaviours in
28 BBR, as well as the interpretation and explanation of the final outcome in human
29 understandable terms. Unfortunately, with a few exceptions, most current EPC
30 approaches are not universal. In BBR, for example, the elementary behaviours
31 are commonly designed off-line for a specific domain or learned from scratch, a
32 significant shortcoming in dynamic or only partially known environments.

33 On the other hand, classical control theory has been explicitly developed to
34 find universal control solutions with an explicit environment model as input [7]. It
35 was also readily extended to hybrid models, where several discrete and continuous
36 components interact in a non-trivial manner (see e.g. overview in [8] and refer-
37 ences therein). In particular, model predictive (or model-following) methods have
38 been found to be applicable to a wide range of control problems and to be efficient
39 at dealing with modelling errors (see e.g. [9, 10]). These methods use a system

¹Notably, BBR is also inherently hybrid, since distinct behaviours can be designed using completely different methodologies: while some of them can use fuzzy logic, others may include a learning algorithm or simply be reactive. However, *EPC* and *hybrid* are, in general, distinct properties.

40 model to generate predictions on the system development, and compute a control
41 signal to optimise this predicted behaviour. Furthermore, the methodology read-
42 ily accepts various learning techniques, both to calculate the control signal and to
43 adaptively calibrate the model in dynamic or partially known environments. How-
44 ever, the detail of the model calibration may vary according to the imposed system
45 structure and dynamics assumptions. For instance, in reinforcement learning ar-
46 chitectures, such as Dyna [11], model corrections are local to the current environ-
47 ment state. Dyna’s principles are also echoed by the modern Bayesian techniques
48 where a POMDP model is recovered while finding the reward maximising policy
49 (e.g. [12, 13]). However, the success of these works has been conditioned on the
50 domain being well factored or on the presence of an oracle to query for the true
51 system state. Furthermore, these approaches can not address the problem of an
52 environment that drifts through a continuous range of models due to their rigid as-
53 sumptions on system structure. To address this issue, much stronger, hybrid control
54 methods have been constructed, usually based on the model predictive (or model-
55 following) principle (see e.g. [14–16]). Some methods even provide theoretical
56 guarantees [14], however at the price of requiring additional modifications to work
57 with discrete space domains or losing this capability entirely.

58 Given these complementary strengths, the fusion of EPC with model-based
59 control can potentially lead to an extremely powerful framework. It would combine
60 the egocentric autonomous representation, i.e. dynamic system without external
61 control input, of a task and the capability to incorporate high level environment
62 knowledge in the form of a system model. Unfortunately, various as they are,
63 classic control theory approaches have an important underlying assumption: the
64 subject of the optimality criteria are the state and the dynamics of the environment.
65 Be that the expected accumulated cost of the state variation (e.g. the classic work of
66 Stengel [7]), be that the proximity to an ideal distribution over system trajectories
67 (e.g. [17]) or be that the cost of system stability (e.g. [18]), the optimality criteria
68 always comes back to consider the underlying system state transitions as the utility
69 source, even if the environment model contains observed quantities only (e.g. [19]).
70 By so doing, this assumption explicitly contradicts the EPC point of view, which
71 hinders the aforementioned fusion of the two control principles.

72 In fact, the only control algorithm that possesses a complete fusion of both the
73 model-based control principles and the EPC view is the Extended Markov Track-
74 ing (EMT) algorithm [20] and its descendants (e.g. [21, 22]). However, as our
75 experiments have revealed, the standard EMT can not cope well with model in-
76 coherences. To this end, in this paper we propose an extended EMT algorithm
77 that has all the aforementioned capabilities: it is an egocentric perceptual control
78 algorithm, it is a universal model-based controller, it is adaptive to environment
79 changes by means of an on-line model calibration, it is a hybrid controller capable

80 of operating in mixed discrete-continuous domains or domains with a hierarchical
81 abstraction of actions. In more detail, for each action available to the agent, we de-
82 ploy an experts ensemble [23] to learn a good estimate of an action’s effects. Such
83 ensembles are known to provide highly flexible and dynamic estimates, which in
84 our case corresponds to fast estimation and calibration of a system model. Notice
85 that this estimate is with respect to the predictive capabilities of the action effects
86 on the agent’s perceptions. Now, the expert ensemble is composed of a finite set
87 of potential effects an action may have, mined from an initial environment model,
88 which are dynamically merged together into a single estimate of an action’s effect.
89 The new control algorithm, the Ensemble Action EMT (EA-EMT) then uses the
90 collection of these estimates to form a complete environment model and proceeds
91 to follow the normal EMT flow of action selection.

92 To demonstrate the adaptive efficacy of the EA-EMT algorithm we have de-
93 vised a set of experiments with various incoherences of the initial system model.
94 In a discrete state environment we have investigated the effects of exogenous per-
95 turbations of three types: *catastrophic* – the environment becomes completely in-
96 consistent with the model, *deviating* – some aspect of the environment behaviour
97 diverges compared to that specified in the model, and *periodic* – the environment
98 alternates between several possible divergences. The results show that EA-EMT re-
99 solved model incoherence and outperformed its EMT predecessor by up to 95%. To
100 clearly demonstrate the hybrid nature and capabilities of the EA-EMT algorithm,
101 we have devised an additional experiment with a continuous state environment,
102 where a task had to be achieved by switching between several pre-specified sub-
103 controllers. In this continuous state environment we have also compared the effects
104 a deviating inconsistency has on EMT-based approaches (both the standard EMT
105 and the EA-EM) and the classical model-following approach. In our experiments,
106 EMT has outperformed the model-following controller under model incoherence,
107 and both have been outstripped by EA-EMT by at least 40% in error rate.

108 To summarise, the contributions of this paper are as follows. First, we intro-
109 duce a new hybrid control method that is equally applicable in environments with
110 discrete, continuous or mixed environment state. This enables the algorithm to
111 serve both as a universal low level mechanism of action selection, and as a high
112 level switching mechanism between separate tuned controllers in a hybrid archi-
113 tecture. In particular, the algorithm is resistant to switching noise, the capability
114 well beyond even the most modern switching methods (e.g. [14]). Second, our ap-
115 proach provides, for first time, an adaptive controller version of the model-based
116 EPC paradigm, enabling in observable terms. Third, EA-EMT is the first algorithm
117 that, without sacrificing its generality with respect to its environment’s continuity,
118 is capable of composing a good control signal even if the underlying environment
119 dynamics are non-stationary, and change over time.

120 The rest of the paper is organised as follows. In Section 2 we detail the operation of the standard EMT Control algorithm. Section 3 follows with the description of our new EA-EMT algorithm, detailing how it reconstructs and calibrates the environment model through the use of expert ensembles. Experimental support for the effectiveness of our approach in handling various model incoherences is given in Sections 4, while the experiments of Section 5 are designed to expose the hybrid nature of our algorithm. To underline the algorithm’s capability to work in environments with changing behavioural trends, our experiments take a special focus on the on-line property of the EA-EMT model calibration. Section 6 summarises the results and gives future directions of this research.

130 2. EMT Control

131 EPC controllers are constructed around some perceptual concept, and necessarily include a subsystem that creates and maintains these perceptions by accumulating and interpreting the observed data. In the case of an EMT Controller the perception is that of the autonomic system dynamics, where the system state appears to stochastically develop over time without external influence. The convenience of this choice is made apparent by the following observation. Assume that some control has been plugged into the environment. The resulting overall system is autonomic, and describes the behaviour of the control-augmented environment in all possible states. Furthermore, although we may not know what specific control law will bring it about, we frequently can describe the autonomic dynamics that we would consider to be ideal or optimal. For example, in IFR, the behaviour of instrument gauge is described without specifying what actions the pilot has to take to achieve this behaviour. This approach is adopted by the EMT controllers, the control task is described by a perception of an idealised autonomic system dynamics, and the algorithm has to sequence actions to achieve the perception of this ideal. To do so, however, the controller requires a subsystem that creates and maintains the necessary perception, and in this paper the subsystem is the Extended Markov Tracking (EMT) algorithm, that also lends its name to the entire control scheme.

149 Formally, the EMT algorithm produces and maintains an estimate of a stochastic state transition function that models the autonomic system behaviour. It does so by performing a conservative update, specifically it minimises the Kullback-Leibler divergence between the new and the old estimate, with the limitation that the new estimate has to match the most recently observed system transition. In more detail, assume that two probability distributions over the system state, p_t and p_{t+1} , are given that describe two consecutive states of knowledge about the system, and τ_t^{EMT} is the old estimate of the system dynamics. Then the EMT update, abbreviated by $\tau_{t+1}^{EMT} = H[p_t \rightarrow p_{t+1}, \tau_t^{EMT}]$, is the solution of the optimisation problem

158 depicted in Fig. 1, where D_{KL} is the Kullback-Leibler divergence. The optimisa-
 159 tion can be recast as finding joint distribution with given marginals, and solved by
 160 an Iterated Proportional Fitting (IPF) procedure [24]. In fact, the EMT update, H ,
 161 can be calculated for any practical set of distributions that describe p_t , p_{t+1} and
 162 τ_t^{EMT} , although in more general situations approximate representations, such as
 163 particle filters or unscented transforms, may be necessary. This enables EMT to
 164 uniformly treat both discrete, continuous and hybrid state spaces. However, to ease
 165 the exposition, in this paper we concentrate on two simplest distribution families,
 166 namely the discrete and the Gaussian distributions. For these families the IPF has
 167 been well studied and needs not to be proof-checked [25, 26].

$$\begin{aligned} \tau_{t+1}^{EMT} &= \arg \min_{\tau} D_{KL}(\tau \times p_t \| \tau_t^{EMT} \times p_t) \\ \text{s.t. } p_{t+1}(x') &= \sum_x (\tau \times p_t)(x', x) \\ \text{and } p_t(x) &= \sum_{x'} (\tau \times p_t)(x', x) \end{aligned}$$

Figure 1: The EMT Update

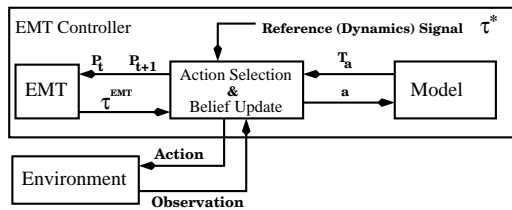


Figure 2: The closed loop of EMT Control

168 To complete the EMT control loop, however, we still need to describe how the
 169 observation data is accumulated to form the perception of the system dynamics. To
 170 this end, we have to address the type of environment models we will be working
 171 with. Although EMT can work with more general environmental descriptions (see
 172 e.g. [22]), it has been more commonly used with a discrete Markovian environ-
 173 ment with partial observability, described by a tuple $MEnv = \langle S, s_0, A, T, O, \Omega \rangle$,
 174 where: S is the set of all possible environment states; $s_0 \in \Delta(S)$ is the initial state
 175 distributions of the environment, where $\Delta(S)$ is a family of distributions over S ; A
 176 is the set of all actions applicable in the environment; $T : S \times A \rightarrow \Delta(S)$ is the
 177 environment's probabilistic transition function, where $T(s'|a, s)$ is the probability
 178 that the environment will move from state s to state s' under action a ; O is the set
 179 of all possible observations; $\Omega : S \times A \times S \rightarrow \Delta(O)$ is the observation probability
 180 function, where $\Omega(o|s', a, s)$ is the probability that o will be observed given that the
 181 environment moved from state s to state s' under action a .

182 This naturally connects with the EMT algorithm, as knowledge about the sys-
 183 tem is summarised by a distribution vector over the system states $p_t \in \Delta(S)$, in
 184 which case the system dynamics estimator created by EMT has the form of a con-
 185 ditional probability $\tau : S \rightarrow \Delta(S)$.

186 Given this, the overall control algorithm, termed *EMT Control*, forms a closed
 187 loop control with a reference signal [7]. Fig. 2 depicts the resulting scheme. Three
 188 sub-modules form the *EMT Controller* that interacts with an *Environment* by ap-

189 plying actions in and receiving observations from it: the *Model*, the *EMT* estimator,
 190 and the decision making module of *Action Selection and Belief Update*. The *Model*
 191 module is queried for the effects T_a of an action a on the real system state. These
 192 effects are used both in predicting future perceptions, and in filtering the observed
 193 data to maintain system state beliefs. The *EMT* module is used to estimate the per-
 194 ceived dynamics τ^{EMT} that explain the change in beliefs about the system from p_t
 195 at time t to p_{t+1} at time $t + 1$. The central, decision making module, interconnects
 196 the *Model* and the *EMT* estimator, and implements the *EMT Control* algorithm,
 197 the detail of which we describe below. Finally, the reference signal, τ^* , encodes
 198 the task to be performed and formally takes the form of the conditional probability
 199 $\tau^* : S \rightarrow \Delta(S)$.

200 Notice that τ^* represents the ideal autonomic system dynamics we would like
 201 to obtain by exercising control. From the EPC point of view, this is the target
 202 perception that we would like to achieve and maintain, hence the standard *EMT*
 203 Control (see Fig. 3) can be described as a greedy one-step look ahead correction
 204 action selection, and it follows a closed loop structure. In more detail, at every
 205 point in time, the algorithm attempts to predict the reaction of an estimation algo-
 206 rithm (*EMT* in this case) to the changes induced by an action (lines 12-16 of the
 207 algorithm), and then chooses the action that shifts the *EMT* estimator closest (line
 208 17) to the reference dynamics τ^* . Once the action has been applied, the response
 209 of the *EMT* estimator to the changes in the environment is registered (line 20), and
 210 the control loops to make its next decision.

211 At this point, we would like to underline the strength of the task representation
 212 by the autonomic system dynamics τ^* . First, while deterministic dynamics are a
 213 way to concisely represent feasible sequences of states, probabilistic dynamics can
 214 also engender a preference over such sequences. Thus, system dynamics τ^* can en-
 215 code a richer variety of preferences and tasks for *EMT* control, than, for instance,
 216 a reward function over states would. Second, in Markov chains, system dynam-
 217 ics completely determine the system state in the long run. As an outcome, the
 218 knowledge about the initial system state is not essential to *EMT* control operation,
 219 expanding its applicability. Furthermore, although over the given state space the
 220 τ^* transition is Markovian, the task it describes needs not be Markovian within the
 221 environment itself. This is due to the fact that the model's state space is abstract,
 222 and each state can serve as a tag for complex, time extended events. Notice, how-
 223 ever, that the controller action selection in lines 12-16 is heavily dependent on the
 224 environment model, as it uses the mapping T_a to predict action effects. However, if
 225 the model is incoherent the reaction of *EMT* can not be estimated correctly, which,
 226 in turn, will lead to selection of a suboptimal action. Thus, in what follows, we
 227 modify the action selection process to vary the environment model it uses.

<p>Require:</p> <p>Set the system state estimator: $p_0(s) = s_0 \in \Delta(S)$</p> <p>Set the system dynamics estimator: $\tau_0^{EMT}(\bar{s} s) = prior(\bar{s} s)$</p> <p>Set time to $t = 0$.</p> <p>11: loop</p> <p>12: for all $a \in \mathcal{A}$ do</p> <p>13: Set $\bar{T}_a = T_a$ {use transition model T directly}</p> <p>14: Set $\bar{p}_{t+1}^a = \bar{T}_a * p_t$</p> <p>15: Set $D_a = H[p_t \rightarrow \bar{p}_{t+1}^a, \tau_t^{EMT}]$</p> <p>16: Set $V(a) = \langle D_{KL}(D_a \tau^*) \rangle_{p_t}$</p> <p>17: Select $a^* = \arg \min_a V(a)$</p> <p>18: Apply a^*, receive observation $o \in \mathcal{O}$</p> <p>19: Compute p_{t+1} due to the Bayesian update: $p_{t+1}(s) \propto \Omega(o s, a) \sum_{s'} \bar{T}(s a, s') p_t(s')$</p> <p>20: Compute $\tau_{t+1}^{EMT} = H[p_t \rightarrow p_{t+1}, \tau_t^{EMT}]$</p> <p>21–25: {no model update}</p> <p>26: Set $t := t + 1$</p>

Figure 3: The standard EMT control algorithm. Note: EA-EMT will modify lines 13,21-25.

228 3. Ensemble Action EMT

229 Although the standard EMT Control is attractive in its combination of the ego-
230 centric control perspective and the task description by the perceived system dy-
231 namics, our experiments (see Section 4) have revealed that its performance de-
232 teriorates significantly if the environment model is incoherent. However, we be-
233 lieve (and will subsequently demonstrate) that, by providing the algorithm with
234 an additional method to correct model incoherences, it is possible to rectify the
235 deterioration. Now, there are many incoherences a Markovian model, $MEnv = \langle$
236 $S, s_0, A, T, O, \Omega \rangle$, may have. Specifically, while the choice of the state, action
237 and observation spaces, as well as the observability function, may be dictated by
238 subjective considerations (e.g. to make it more readable for the human domain
239 designers), the transition function T is always dictated by the environment. Thus,
240 in this work we choose to concentrate on the quality of the transition function T .
241 This function maps actions into stochastic matrices, so that for each action $a \in A$
242 the matrix $T_a = T(\cdot, a)$ models the effects of that action on the system state.
243 The difference between the matrix T_a and the true effects of the action $a \in A$ is
244 the incoherence type we have resolved in the EA-EMT algorithm (Fig. 4). Thus,
245 while the standard EMT Control views the transition mapping, $a \mapsto T_a$, to be con-
246 stant, the EA-EMT algorithm modifies its transition mapping over time, reducing
247 the mapping’s incoherence. However, before we go into the details of how it was
248 implemented, we need to explain the principles of the approach taken by EA-EMT.

249 EA-EMT assumes that, although the mapping $T : A \rightarrow \Delta(S)^S$ is incoherent,
250 the set of matrices $T_A = \{T_a = T(\cdot, a)\}_{a \in A}$ represents feasible effects that the
251 actions may have. The algorithm then attempts to assemble a better mapping,

252 $\bar{T} : A \rightarrow \Delta(S)^S$, based on the set T_A . More specifically, for each action $a \in A$
 253 the transition matrix \bar{T}_a is a weighted linear combination of matrices in the set
 254 T_A , that is $\bar{T}_a = \sum_{b \in A} T_b * w_a(b)$. Intuitively, the weight $w_a(b)$ represents the
 255 similarity between the matrix $T_b \in T_A$ and the effects that the action $a \in A$ has on
 256 on the environment state. As the interaction between the EA-EMT algorithm and
 257 the environment progresses, the weights $w_a(\cdot)$ are updated, modifying the mapping
 258 $\bar{T} : A \rightarrow \Delta(S)^S$ to reduce its incoherence with the environment.

259 The intuition behind this approach stems from *Polytopic Linear Models (PLM)*
 260 with continuous state, where a complex non-linear system is represented as a com-
 261 bination of a finite set of simpler linear sub-systems [15]. Similarly, in our for-
 262 malism, an action $a \in A$ may be more than a primitive operation. Rather, it may
 263 represent a subsystem with a complex underlying controller, which forces the sys-
 264 tem to follow dynamics described by T_a . In fact, by enriching the set T_A , one can
 265 guarantee that environment incoherences of interest will be well captured. As an
 266 utterly extreme example consider a dynamic system with a discrete state space.
 267 By setting T_A to be the set of permutation matrices, we essentially allow \bar{T}_a to be
 268 any matrix from the polytope of stochastic matrices, and endow EA-EMT with the
 269 capability to capture any environment disturbance, be it a randomly reoccurring
 270 one or be it a disturbance localised to a particular system state. Although the re-
 271 lationship between the composition of T_A and its expressiveness needs not be this
 272 extreme, and in practice only small sized T_A is required, its exact properties are
 273 non-trivial. In fact, it forms a separate branch of research, where the works by An-
 274 gelis [15] and Cesa-Bianchi [23] are only few representatives of a vast literature,
 275 that falls out of scope of this paper. Nevertheless, we can safely assume that T_A
 276 forms a sufficiently large polytope that includes all relevant system dynamics.

277 Now, the update of the weights $w_a(\cdot)$ is based on the approach of predictions
 278 with expert ensembles [23]. The intuition behind this approach is that, when mak-
 279 ing a prediction or a decision, a readily available set of feasible alternatives (the
 280 expert ensemble) can be merged together to form a prediction which is potentially
 281 better than any of the alternatives standing alone. The dynamic properties of this
 282 merger are such, that it can be readily applied even if the best prediction (or the
 283 best decision) is not stationary, but rather changes over time. This made the choice
 284 of expert ensembles particularly attractive to maintain a system model in varying,
 285 unstable environments. Specifically, in our algorithm the *expert ensemble* is the set
 286 T_A , where each expert attempts to predict the effects an action would have on the
 287 environment state. From this point of view, the weight $w_a(b)$ expresses how much
 288 the expert $T_b \in T_A$ is trusted to capture the effects of the action $a \in A$ correctly.
 289 Once EA-EMT has applied an action, a^* , it measures the discrepancy between the
 290 effect a^* had and the effect predicted by expert T_b . The lower the discrepancy, the
 291 higher will be the weight $w_{a^*}(b)$ when the next control decision is made.

292 Given the above principles, we have modified the standard controller algo-
 293 rithm. Specifically, line 13, previously directly substituted into the calculations the
 294 transition function from the provided model. Whereas now it uses a weighted com-
 295 bination of the matrices in T_A , which is continually tuned by the expert ensemble
 296 to improve its representation of an action's effects. The rest of the computations
 297 proceed as before until the EMT estimate, τ_{t+1}^{EMT} , of the action outcome is com-
 298 puted in line 20: the algorithm predicts the effects of each action on the EMT
 299 estimate, chooses the action that would bring τ_{t+1}^{EMT} closest to the reference signal
 300 τ^* , applies the action and receives an observation. At that point, the algorithm has
 301 to measure the performance of each expert, and update the weights. Now, recall
 302 that the algorithm operates in terms of subjective beliefs, the relevant effects of
 303 the action are thus those expressed in the EMT estimate τ_{t+1}^{EMT} . This means that
 304 the performance of each expert can be expressed by the distance between the es-
 305 timate τ_{t+1}^{EMT} and the estimate that would have been obtained based on the expert
 306 prediction. This distance is computed in lines 22-24, and the weight of the expert is
 307 updated accordingly. Specifically, the old weight of the expert is multiplied by β^d ,
 308 where $\beta \in (0, 1)$ is the parameter of the update and d is the distance above. Once
 309 all weights are updated, they are normalised to sum to 1, so that \tilde{T}_a at the next step
 310 will be a stochastic matrix. Notice that all these operations take time polynomial in
 311 the model parameters, such as the size of state, action and observation spaces. This
 312 makes EA-EMT a computationally efficient and scalable algorithm, an attractive
 313 property systems where environment models tend to be large.

<p>Require: ... Set action weight vectors: $w_a(a') \propto \delta_a(a') + \epsilon$ Set time to $t = 0$.</p> <p>11: loop 13: Set $\tilde{T}_a = \sum_{a'} T_{a'} * w_a(a')$ 21: for all $a \in \mathcal{A}$ do 22: Set $\tilde{p}_{t+1}^a = \tilde{T}_a * p_t$ 23: Set $D_a = H[p_t \rightarrow \tilde{p}_{t+1}^a, \tau_t^{EMT}]$ 24: Set $V(a) = \langle D_{KL}(D_a \tau_{t+1}^{EMT}) \rangle_{p_t}$ 25: Set $w_{a'}(a) \propto w_{a'}(a) \beta^{V(a)}$ 26: Set $t := t + 1$</p>
--

Figure 4: The EA-EMT control algorithm: only changes to the standard EMT control are shown.

314 4. Experimental Evaluation: Discrete State Space

315 To test the effectiveness of the EA-EMT algorithm, we have devised a set of com-
 316 parative tests with the standard EMT Controller. The latter is a natural baseline,

317 as it is the only other universal control algorithm capable of complete fusion of
 318 the EPC and the model-based paradigms. In discrete state systems this is also the
 319 only baseline, as no other control algorithm can reproduce the action selection se-
 320 quence of EMT. Fortunately, in the environments with a continuous state space,
 321 which we tend to in Section 5, the sequence of actions selected by EMT can be at
 322 least partially reproduced by model-following control algorithms, and we imme-
 323 diately use it to provide an additional baseline comparison. In all cases, we have
 324 preferred a simulated system so that the true effects of our control algorithm will
 325 not be confused with the properties of an embodied physical system.

326 Now, to support comparability with previous work on EMT variations, all tests
 327 were based on modifications of the Drunk Man (D-Man) domain: a controlled
 328 random walk over a linear graph (see Fig. 5 for the principle structure) with ac-
 329 tions weakly modulating the probability (only a small discrete set of probabilities
 330 in the range $(\epsilon, 1 - \epsilon)$ with $\epsilon \gg 0$ is attainable) of the left and the right steps.
 331 The domain is also partially observable, namely, instead of its true position on the
 332 graph, an agent receives as an observation a random position within the two step
 333 neighbourhood of agent's location. In turn, a task within the domain is represented
 334 by a conditional probability $\tau^*(s'|s)$, the reference signal for the controller, spec-
 335 ifying what sort of motion through the state space has to be induced. During an
 336 experiment run, the control algorithm was provided with a Markovian environment
 337 model, $MEnv = \langle S, s_0, A, T, O, \Omega \rangle$, incoherent with the true behaviour of the do-
 338 main. The incoherences were created by introducing exogenous perturbations to
 339 the behaviour of the D-Man domain. In particular, three perturbations, making the
 340 model of the standard D-Man domain increasingly incoherent with the actual envi-
 341 ronment behaviour, were used: **Deviating**, where an additional deterministic step
 342 (to the right) was done; **Periodic**, where the direction of an additional deterministic
 343 step changed over time; and **Catastrophic**, where a random permutation of actions
 344 was selected $\sigma : A \rightarrow A$, so that when the controller applied action $a \in A$, the
 environment responded instead to $\sigma(a)$. Three baselines were obtained in various

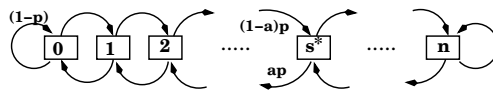


Figure 5: Principle structure of the Drunk Man domain.

345 combinations: standard EMT Control algorithm operating in a perturbed environ-
 346 nment, standard EMT Control operating within an unperturbed environment, and
 347 standard EMT Control operating in a perturbed environment with its model cor-
 348 rectly encoding the environment perturbation. At least two baselines are present
 349 in each experimental setting to provide comparative performance bounds and the
 350

351 99.5% confidence envelope is depicted in all plots. In all our experiments the ref-
 352 erence dynamics for the controller is given by $\tau^*(s'|s) \propto \delta_{s^*}(s') + \epsilon$, where $\epsilon > 0$
 353 is small. In other words, the target prescribes that the environment should almost
 354 surely move to the ideal state s^* from any other state. In our experiments the state
 355 space was $S = \{0, \dots, 12\}$, and the ideal state $s^* = 6$. Notice that, due to the prob-
 356 abilistic nature of the domain, any reasonable² control scheme set to accomplish
 357 the task would result in a bell shaped empirical distribution of the system state.
 358 Success of the control scheme can then be readily appreciated visually by the dif-
 359 ference of the expected value and the ideal system state, as well as the standard
 360 deviation of the empirical state distribution. The empirical distribution was taken
 361 over a 200 decision step *sliding window*, to obtain statistically significant distri-
 362 bution shape. In turn, the overall length of experimental runs was then chosen
 363 to be sufficiently large to enable analysis of stable trends of the empirical 200-
 364 step distribution. In particular, for the *catastrophic* and the *deviating* perturbations
 365 each experiment run was 1000 steps. The necessity to obtain statistical signifi-
 366 cance while preventing the algorithm from completely stabilising, has also led to
 367 the choice of the 500 step period for the *periodic* perturbation experiments, accom-
 368 panied by the 5000 step total length of each experiment run. Although alternative
 369 experimental setups were also run, varying both the sliding window size and the
 370 experiment length, their results were similar, we, therefore, omit them due to space
 371 limitations. Nevertheless, the aforementioned sequence of choices is reflected in
 372 the way our experimental results are presented: *deviating*, *catastrophic* and then
 373 *periodic* perturbations. Furthermore, to present an overall evaluation of a control
 374 scheme's performance, rather than a comparison of multiple parameters, we also
 375 measured the distance between the empirical distribution and δ_{s^*} using l_1 norm.

376 To further the intuition of this domain, consider once more the IFR example
 377 where the pilot has to maintain flight level within the air corridor prescribed by the
 378 ground control. If we discretise the space of possible flight levels we will obtain a
 379 linear graph depicted in Fig. 5. The transitions between the states are controlled,
 380 but are also subject to random changes in the air density or wind gusts. Ideally, the
 381 auto-pilot will need to actively return the airplane to the ideal, centre flight level.

382 4.1. Deviating Perturbation

383 In this experiment we introduce a deviating perturbation. That is, beyond the usual
 384 probabilistic step, the environment has also deterministically shifted in one direc-
 385 tion along the linear graph. For example (referring to Fig. 5) if the system reached

²Unreasonable, for instance, would be choosing a constant action to equalise the left and the right step probabilities, as this would result in an almost uniform distribution, utterly defeating the controller purpose.

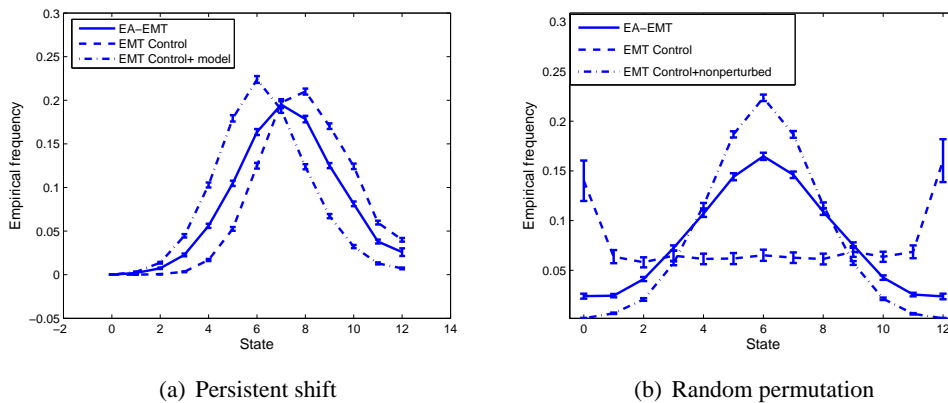
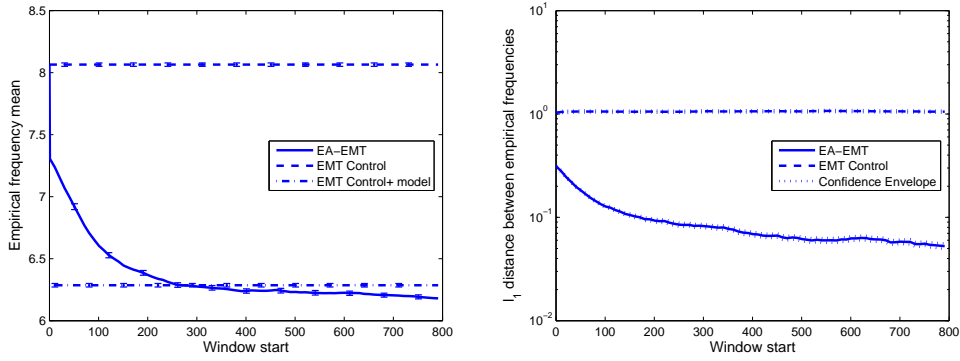


Figure 6: EA-EMT performance under (a) deviating and (b) catastrophic perturbations

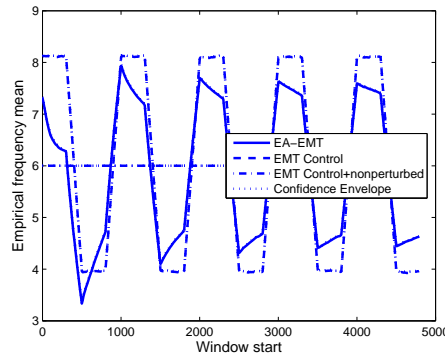
386 state $k \in \{0, \dots, n - 1\}$, the additional step will shift it to state $k + 1$. In this con-
 387 text, Fig. 6(a) shows the empirical distribution of system states under three control
 388 strategies: the EA-EMT controller and the standard EMT Controller equipped with
 389 the standard D-Man model (thus excluding the shift modelling), and the standard
 390 EMT Controller equipped with the environment model that explicitly captures the
 391 additional shift. The figure shows the complete empirical distribution of the EA-
 392 EMT obtained during the first 200 control choices made in this experiment, and
 393 marks a definitive improvement in performance. This can be seen from the fact
 394 that the standard EMT Control fails to enforce the reference dynamics τ^* , with
 395 the system spending the majority of its time away from the ideal state, $s^* = 6$,
 396 while EA-EMT manages to force the state distribution to concentrate closer to s^* .
 397 In fact, the distance between δ_{s^*} and the EA-EMT distribution induced in the first
 398 200 steps is 40% less than the comparable distance for the EMT controller. This,
 399 however, does not fully reflect the adaptability of EA-EMT. To this end, Fig. 7(a)
 400 shows how the mean of the empirical distributions of the 200 step windows behave.
 401 The distributions induced by EMT Control do not change over time, resulting in
 402 straight horizontal lines depicting the constancy of the mean. On the other hand,
 403 the data shows that EA-EMT quickly adapts, the algorithm induces the empirical
 404 state distribution with the mean approaching the ideal state $s^* = 6$. In this respect,
 405 EA-EMT even slightly surpasses the performance of the standard EMT algorithm
 406 with the correct environment model. This is due to the adaptive portion of EA-
 407 EMT contributing to the tie breaking when considering similar actions – this tie
 408 breaking is rigid in EMT Control. Similar pictures occur with respect to the vari-
 409 ance of the empirical distributions. This means that EA-EMT overcomes the model

410 incoherence and increasingly concentrates the state empirical distribution around
 411 the ideal state, which is exactly what the reference dynamics, τ^* , requires.



(a) Persistent Shift

(b) Random Permutation



(c) Switching Shift

Figure 7: EA-EMT adaptation to various perturbations. Notice the log-scale of the Y axis in (b).

411

412 **4.2. Catastrophic Perturbation**

413 The action space of the D-Man domain has a simple intuitive interpretation – the
 414 action sets how quickly the system state will shift left or right. The deviating per-
 415 turbation did not exceed this interpretation, it simply meant that the system will
 416 naturally move in one direction faster than the other. In a way it also meant that
 417 the perturbation induced a very mild model incoherence – principally the model re-
 418 mained correct. However, EA-EMT can adapt to much more severe model incoher-
 419 ences. In fact, in the next set of experiments the environment model is completely
 420 incorrect. For each run in this experiment set a random permutation $\sigma : A \rightarrow A$

421 was selected. Then, when action $a \in A$ was applied, the environment reacted as if
422 the action was $\sigma(a)$.

423 In more depth, Fig. 6(b) shows the empirical distributions obtained in the first
424 200 steps of decision making. Permuting the action breaks any connection between
425 what EMT Control expects the action to do and what actually occurs in the environ-
426 ment, essentially the actions are scrambled and the EMT Control chooses a random
427 action. This results in the algorithm’s failure – the empirical state distribution is
428 equivalent to that of applying no control at all, with higher probability of terminal
429 states due to the failure of the respective left and right steps. In contrast, EA-EMT
430 easily adapts and performs increasingly well, as can be seen in Fig. 7(b). Follow-
431 ing the development of the empirical distribution within a 200 step sliding window,
432 the figure shows the l_1 distance from the distribution formed by the standard EMT
433 algorithm in the non-perturbed environment. This data demonstrates that EA-EMT
434 exponentially quickly discovers the true effects of actions and approaches the per-
435 formance of the EMT control in a non-perturbed environment. Even though the
436 empirical distribution of the first 200 steps includes the first decisions made based
437 on the scrambled model, it already recovers 70% of the performance lost due to the
438 model incoherence and, through further adaptation, it reaches 95% recovery.

439 4.3. Periodic Perturbation

440 Finally, it is important to ensure that the algorithm can perform well in a dynami-
441 cally changing environment. For example, a robot’s body is subject to environmen-
442 tal effects, and its response to control will change accordingly. Some environment
443 parameters, like the daily temperature variation on Lunar surface, may be extreme
444 and persistently reoccurring. To test EA-EMT in such environments, we consider
445 yet another perturbation: an additional deterministic step is made, and the direction
446 of the step switches between left and right with constant period (500 control steps
447 in our experiments). The shape of the distributions formed by the controllers are
448 equivalent to those in the persistent shift experiment (see Fig. 6(a)), and we omit
449 the respective graph. On the other hand, the development of the empirical distri-
450 bution over time is quite different. In particular, Fig. 7(c) shows the behaviour of
451 the mean value for empirical distributions calculated within a 200 step sliding win-
452 dow. While the standard algorithm literally switches from one value to another,
453 depending on the direction of the shift, the performance of EA-EMT always shows
454 recovery after a direction switch occurs. Notice also, that the magnitude of the
455 mean variation at the switch point becomes significantly (25%) less for EA-EMT
456 than the standard EMT. This suggests that, beyond its ability to recover from ir-
457 relevant adaptations, the adaptive controller version learns to reduce the control
458 inertia. In other words the algorithm reduces the impact of the sudden change in
459 the environment behaviour, stabilising the overall performance.

460 **5. Experimental Evaluation: Continuous State Space**

461 To complete the demonstration of our algorithm, we apply EA-EMT to a continu-
 462 ous state environment, where a task is achieved by switching between pre-specified
 463 sub-controllers. This combination of the discrete switching and the continuous
 464 switching components, clearly show EA-EMT to be a hybrid controller. Notably,
 465 neither the structure nor the principle of application change with the transition from
 466 a discrete to a continuous state space. The transition is achieved simply by replac-
 467 ing the finite dimensional vector, that has represented state probability distribution
 468 in the discrete case, by a Gaussian distribution to represent the state distribution
 469 of the continuous domain. Similarly, stochastic transition matrix is replaced by a
 470 conditional Gaussian distribution to capture system dynamics. Furthermore, the
 471 amount of underlying calculations grows only polynomially with the dimension of
 472 the state and the observation spaces. It is this computational scalability, and the fact
 473 that no modification is required nor made to the reasoning of the action selection
 474 procedure, which remains fully and completely intact whatever the environment
 475 dimensionality is, that grant EA-EMT almost universal applicability. It allows our
 476 algorithm to be deployed both as a direct low level controller, and as a part of a
 complex hierarchical hybrid controller with multiple levels of abstraction.

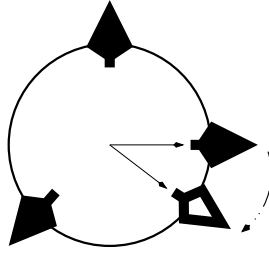


Figure 8: Hovercraft scheme.

477
 478 The specific domain we chose is that of a hovercraft with three thrusters de-
 479 picted in Fig. 8. Solid arrows show thruster directions, while the hollow arrow
 480 denotes a potential mistake in that thruster's model. From the perspective of our
 481 IFR example, such modelling mistake would correspond to a sudden change in the
 482 plane's responses, for instance due to a collision with a bird or a mechanical mal-
 483 function. The system generically develops in discrete time using the equation:

$$\begin{bmatrix} x_{k+1} \\ \dot{x}_{k+1} \\ y_{k+1} \\ \dot{y}_{k+1} \end{bmatrix} = \begin{bmatrix} 1 & h & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & h \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_k \\ \dot{x}_k \\ y_k \\ \dot{y}_k \end{bmatrix} + \begin{bmatrix} \frac{h^2}{2} & 0 \\ h & 0 \\ 0 & \frac{h^2}{2} \\ 0 & h \end{bmatrix} [v_1, v_2, v_3] \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix}$$

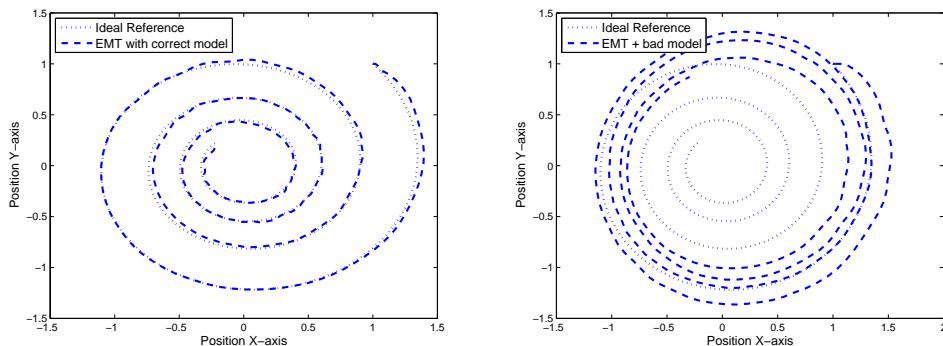
485 In the equation, v_i denotes the directional force distribution of a thruster, u_i

486 its basic level of activity, and h denotes the time span during which the thrust
487 was applied. To further underline the use of EA-EMT as a switching mecha-
488 nism of a hybrid controller, we restrict u_i in our experiments to a finite discrete
489 set. Specifically, we used 5 activation levels between 0.2 and 1.0 in equal inter-
490 vals, and only one thruster could have a non-zero activation at any time, so that
491 total of 15 distinct joint activations were possible. This naturally simulates the
492 situation that occurs in hybrid systems, where an action corresponds to the ap-
493 plication of a distinct sub-system controller, rather than a choice of a continuous
494 value control signal. Two configurations of thrusters were used. Configuration A,
495 $[v_1, v_2, v_3] = \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & -1 \end{bmatrix}$, that corresponds to the solid arrows in Fig. 8; and config-
496 uration B, $[v_1, v_2, v_3] = \begin{bmatrix} 1 & 0 & -1 \\ -0.3 & 1 & -1 \end{bmatrix}$, that corresponds to a structural failure of a
497 thruster depicted by the hollow arrow. In all experiments, while the controller al-
498 gorithm was given either the environment model with thruster Configuration A or
499 B, the actual motion of the hover craft was always simulated using Configuration
500 A. This discrepancy allowed us to test the performance of our algorithm under a
501 deviating modelling incoherence.

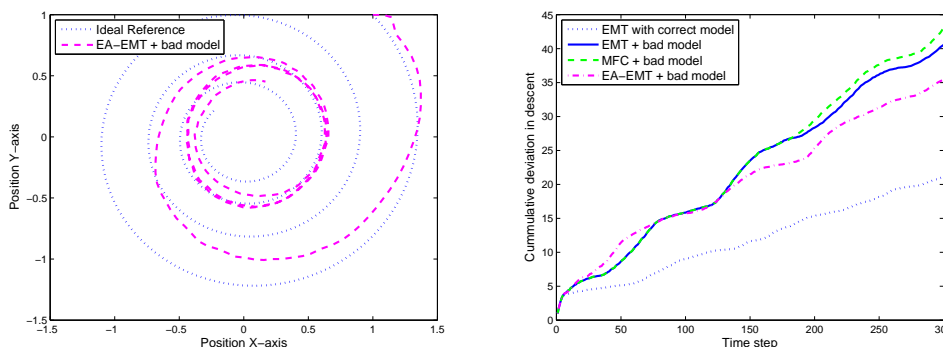
502 Now, to provide a quantitative performance measure, we have set several con-
503 trol algorithms with the task to simulate a gradual spiralling descent towards zero
504 from rest at coordinates $[1, 1]$, which we have described by an autonomic linear
505 system with the equation given below. Recalling once more our IFR scenario, such
506 system would correspond, for example, to the necessary relative properties of the
507 altitude and speed of the airplane, as well as their development in time, during a
508 landing procedure. As before, h denotes the time span of a single step, while λ de-
509 notes the decay of the spiral and θ the rotation angle of a single step of the system.

$$510 \begin{bmatrix} x_{k+1} \\ \dot{x}_{k+1} \\ y_{k+1} \\ \dot{y}_{k+1} \end{bmatrix} = \begin{bmatrix} \lambda \cos(\theta) & 0 & -\lambda \sin(\theta) & 0 \\ \frac{2}{h}(\lambda \cos(\theta) - 1) & -1 & -\frac{2}{h}\lambda \sin(\theta) & 0 \\ \lambda \sin(\theta) & 0 & \lambda \cos(\theta) & 0 \\ \frac{2}{h}\lambda \sin(\theta) & 0 & \frac{2}{h}(\lambda \cos(\theta) - 1) & -1 \end{bmatrix} \begin{bmatrix} x_k \\ \dot{x}_k \\ y_k \\ \dot{y}_k \end{bmatrix}$$

512 In more detail, the algorithms we have considered were EMT, EA-EMT and
513 a discrete Model Follower Controller (MFC). The latter algorithm has been se-
514 lected for its robustness and ubiquity of its principle (see e.g. [7, 14, 16]), making
515 it suitable to produce a baseline comparison. The MFC algorithm operated in the
516 usual manner, specifically, given the current hover coordinates, the algorithm se-
517 lected thrust to minimise the discrepancy between the outcome predicted by the
518 task's equation and the equation of the hovercraft's model. We have tuned EMT
519 initialisation and task representation parameters so that, for the Configuration A
520 thrusters model, its decisions coincide with MFC. We conjecture, in fact, that EMT
521 is formally a more general approach than MFC, in that EMT can always be tuned
522 to reproduce MFC's behaviour. The resulting hovercraft trajectory is depicted in



(a) EMT (and MFC) with Configuration A (correct) model (b) EMT with Configuration B (wrong) model



(c) EA-EMT with Configuration B model (d) Cumulative error of controllers

Figure 9: Hovercraft trajectories under various controllers algorithms and controller models and their cumulative error. In all cases physical simulation adopts thruster Configuration A.

523 Fig. 9(a). The dotted line represents the ideal trajectory that could have been ob-
 524 tained if the thrusts u_i were continuous, rather than discretised. The Fig. 9(a)
 525 also demonstrates that the task we posed can be indeed solved by an application
 526 of the standard EMT controller or MFC, forming a performance baseline where
 527 the environment develops exactly as the controller’s model describes it. In turn
 528 Fig 9(b) and Fig 9(c) depict the performance of the EMT and EA-EMT algorithms
 529 provided with Configuration B (wrong) thruster model. Due to the aforementioned
 530 EMT tuning, even under Configuration B the trajectories of EMT and MFC are
 531 extremely similar, and we omit the latter due to space limitations.

532 However, Fig.s 9(a), 9(b), 9(c) can only provide an intuition as to how various
 533 algorithms cope with the task. To clearly distinguish and evaluate the control algo-

534 rithms’ performance we have calculated the cumulative error of these trajectories.
535 That is, for each experiment run at each time step we have computed the difference
536 between the system state that results from the discrete level of thrust chosen by
537 a control algorithm and the system state that resulted from the application of the
538 analytically computed continuous thrust. Fig. 9(d) depicts the accumulation of that
539 discrepancy over time. Initially slightly worse, due to slack expert ensemble initial-
540 isation, over time EA-EMT significantly outperforms both EMT and MFC. Perhaps
541 to further underline the strength of the EMT-based approach in general, notice that,
542 under model incoherence, even the standard EMT outperforms MFC, and aggre-
543 gates trajectory error at a lower rate. Notice that due to thrust discretisation zero
544 error is unachievable, as is witnessed by the error accumulation of EMT (and MFC
545 since they coincide in this case) with the correct Configuration A thrusters model.
546 Furthermore, we have calculated the accumulated thrust utilised by all algorithmic
547 solutions when faced with the bad Configuration B model. The results are given in
548 Table 1. The data confirms that EA-EMT recovers significant portion of losses due
549 to model incoherence. Furthermore, to complete our investigation, we have also
550 measured the amount of energy consumed by the control algorithms in terms of
551 the applied thrust vector norm (see the third column of Table 1). Although at first
552 sight it may look that EA-EMT has conserved some energy by a faster move to a
553 lower spiral loop, in fact, and unlike a passive descent under a gravitational pull,
554 maintaining a tighter trajectory at the same speed necessitates ever higher energy
555 levels to counter the centrifugal force. We are, therefore, inclined to conclude that
556 the energy conservation is an algorithmic property of EA-EMT.

Algorithm/Thruster Configuration	Total Energy	Total Trajectory Discrepancy
EMT(MFC)/Configuration A	132.6	21.067285
MFC/Configuration B	150.8	43.036976
EMT/Configuration B	140.4	40.56558
EA-EMT/Configuration B	117.6	35.400698

Table 1: Total trajectory discrepancy and energy consumption over 300 steps

557 6. Conclusions and Future Work

558 In this paper we present the Ensemble Action EMT algorithm – a control solution
559 that has three important properties: it is an *egocentric perceptual controller*; it is
560 a *universal model-based controller*; it is an *on-line model calibrating controller*;
561 and it is a *hybrid controller* capable of operating in mixed discrete-continuous or
562 hierarchical action abstraction domains. As an EPC solution, EA-EMT describes
563 the control task and the optimality criteria in terms of the agent’s interpretation

564 of sensory input, thus enabling an autonomous agent to formulate internal control
565 tasks, rather than just following an external command. Being a universal model-
566 based solution, EA-EMT is capable of utilising a given environment model, but is
567 not bound to one model or one environment in particular. Finally, on-line model
568 calibration enables EA-EMT application to changing or simply poorly modelled
569 environments.

570 EA-EMT is unlike other adaptive control algorithms based on expert ensem-
571 bles, where experts directly produce actions or plans to be fused (e.g. [16, 27,
572 28])³. Rather, EA-EMT operates in two distinct modules: the expert-based model
573 estimation and a control algorithm that utilises that model. This enables greater
574 design flexibility, and generalisation, particularly with respect to the model type
575 that experts produce. For instance, in robotic soccer – a domain well known to
576 attract hybrid control solutions – environment models are frequently found at the
577 edge of logic and probability-based approaches, especially in opponent plan recog-
578 nition [29–31]. Nevertheless, because of the employed probabilistic notions, these
579 models can still be successfully weighted and fused, albeit necessitating an infer-
580 ence process to do so [29–31]. Furthermore, they still can be evaluated and com-
581 pared via the Kullback-Leibler divergence. As a result, EA-EMT can be expanded
582 to operate even in such a highly complex and dynamic environment as robotic
583 soccer. In fact, the on-line adaptability of the EA-EMT and its computational effi-
584 ciency will be particularly useful.

585 Finally, we also would like to investigate the possibility of altering the weight
586 adaptation to include *forgetting* (inherent tendency of weights to equalise over
587 time) and *update extrapolation* (simultaneous weight modification of actions with
588 similar effects). In particular, forgetting and update extrapolation can serve well
589 in combination with learning approaches. Specifically, we would like to consider
590 the situation where a library of behaviour primitives (or experts) is dynamically
591 composed (see e.g. MOSAIC [16]). In this case, the appearance of new control
592 sub-systems can be handled better, if the expert mixture can be initialised, rather
593 than learned over time, by means of *update extrapolation*. Similarly, older sub-
594 systems can be phased out more effectively if *forgetting* is applied.

595 [1] W. T. Powers, Behavior: The control of perception, Aldine de Gruyter, 1973.

596 [2] S. Thrun, Bayesian landmark learning for mobile robot localization, Machine
597 Learning 33 (1) (1998) 41–76.

³Notably, these methods, particularly Haruno et al. [16], naturally assume control signal metric, which we do not, therefore allowing for more abstract action spaces

- 598 [3] A. Lazanas, J. Claude Latombe, Landmark-based robot navigation, in: Algo-
599 rithmica, 1992, pp. 816–822.
- 600 [4] R. C. Arkin, Behavior-Based Robotics, MIT Press, 1998.
- 601 [5] M. M. Taylor, Editorial: Perceptual control theory and its application, Inter-
602 national Journal of Human-Computer Studies 50 (6) (1999) 433–444.
- 603 [6] W. T. Bourbon, Perceptual control theory, in: H. L. Roitblat, J.-A. Meyer
604 (Eds.), Comparative approaches to cognitive science, MIT Press, 1995.
- 605 [7] R. F. Stengel, Optimal Control and Estimation, Dover Publications, 1994.
- 606 [8] Z. Sun, S. S. Ge, Analysis and synthesis of switched linear control systems,
607 Automatica 41 (2005) 181–195.
- 608 [9] M. Morari, J. Lee, Model predictive control: Past, present and future, Com-
609 puters and Chemical Engineering 23 (9) (1999) 667–682.
- 610 [10] J. Morningred, B. Paden, D. Seborg, D. Mellichamp, An adaptive nonlinear
611 predictive controller, Chem. Eng. Sci 47 (4) (1992) 755–765.
- 612 [11] R. S. Sutton, Integrated architectures for learning, planning, and react-
613 ing based on approximating dynamic programming, in: Proceedings of the
614 ICML, 1990, pp. 216–224.
- 615 [12] P. Poupart, N. Vlassis, Model-based bayesian reinforcement learning in par-
616 tially observable domains, in: Proceedings of the ISAIM, 2008.
- 617 [13] R. Jaulmes, J. Pineau, D. Precup, A formal framework for robot learning and
618 control under model uncertainty, in: IEEE ICRA, 2007.
- 619 [14] L. Giovanini, A. W. Ordys, M. J. Grimble, Adaptive predictive control us-
620 ing multiple models, switching and tuning, International Journal of Control,
621 Automation, and Systems 4 (6) (2006) 669–681.
- 622 [15] G. Angelis, System analysis, modelling and control with polytopic linear
623 models, Ph.D. thesis, University of Eindhoven (2001).
- 624 [16] M. Haruno, D. M. Wolpert, M. Kawato, MOSAIC model for sensorimotor
625 learning and control, Neural Computation 13 (2001) 2201–2220.
- 626 [17] M. Karny, T. V. Guy, Fully probabilistic control design, Systems and Control
627 Letters 55 (4) (2006) 259–265.

- 628 [18] A. Robertsson, On observer-based control of non-linear systems, Ph.D. thesis,
629 Department of Automatic Control, Lund Institute of Technology (1999).
- 630 [19] M. R. James, S. Singh, M. L. Littman, Planning with predictive state repre-
631 sentations, in: Proceedings of the ICMLA, 2004, pp. 304–311.
- 632 [20] Z. Rabinovich, J. S. Rosenschein, Extended Markov Tracking with an appli-
633 cation to control, in: The 1st MOO Workshop, 2004, pp. 95–100.
- 634 [21] Z. Rabinovich, J. S. Rosenschein, Multiagent coordination by Extended
635 Markov Tracking, in: The 4th AAMAS, 2005, pp. 431–438.
- 636 [22] A. Adam, Z. Rabinovich, J. S. Rosenschein, Dynamics based control with
637 PSRs, in: 7th AAMAS, 2008, pp. 387–394.
- 638 [23] N. Cesa-Bianchi, G. Lugosi, Prediction, learning, and games, Cambridge
639 University Press, 2006.
- 640 [24] S. Kullback, Probability densities with given marginals, *The Annals of Math-*
641 *ematical Statistics* 39 (4) (1968) 1236–1243.
- 642 [25] E. Cramer, Conditional iterative proportional fitting for Gaussian distribu-
643 tions, *Journal of Multivariate Analysis* 65 (2) (1998) 261–276.
- 644 [26] S.-C. Fang, J. R. Rajasekera, H. S. J. Tsao, *Entropy Optimization and Math-*
645 *ematical Programming*, Kluwer Academic Publishers, 1997.
- 646 [27] E. Even-Dar, S. M. Kakade, Y. Mansour, Experts in a Markov decision pro-
647 cess, in: NIPS, 2004.
- 648 [28] B. Argall, B. Browning, M. Veloso, Learning to select state machines using
649 expert advice on an autonomous robot, in: ICRA, 2007.
- 650 [29] H. Bui, A general model for online probabilistic plan recognition, in: 18th
651 IJCAI, 2003, pp. 1309–1315.
- 652 [30] P. Riley, M. Veloso, Recognizing probabilistic opponent movement models,
653 in: The 5th RoboCup Competitions and Conferences, 2002.
- 654 [31] D. V. Pynadath, M. P. Wellman, Probabilistic state-dependent grammars for
655 plan recognition, in: 16th UAI, 2000, pp. 507–514.