Deep Lagrangian Networks for end-to-end learning of energy-based control for under-actuated systems

Michael Lutter¹, Kim Listmann² and Jan Peters^{1,3}

Abstract— Applying Deep Learning to control has a lot of potential for enabling the intelligent design of robot control laws. Unfortunately common deep learning approaches to control, such as deep reinforcement learning, require an unrealistic amount of interaction with the real system, do not yield any performance guarantees, and do not make good use of extensive insights from control theory. In particular, common black-box approaches – that abandon all insight from control – are not suitable for complex robot systems.

We propose a deep control approach as a bridge between the solid theoretical foundations of energy-based control and the flexibility of deep learning. To accomplish this goal, we extend Deep Lagrangian Networks (DeLaN) to not only adhere to Lagrangian Mechanics but also ensure conservation of energy and passivity of the learned representation. This novel extension is embedded within a energy control law to control underactuated systems. The resulting DeLaN for energy control (DeLaN 4EC) is the first model learning approach using generic function approximation that is capable of learning energy control because existing approaches cannot learn the system energies directly. DeLaN 4EC exhibits excellent real-time control on the physical Furuta pendulum and learns to swing-up the pendulum while the control law using system identification does not.

I. INTRODUCTION

Control laws are essential to achieve intelligent robots that enable industrial automation, human-robot interaction and locomotion. The common approach is to manually derive the system dynamics, measure the masses, lengths, inertias of the disassembled mechanical system [1] and finally use these equations to engineer a control law for this specific system. Therefore, this engineering approach requires significant effort. In stark contrast, many learning to control approaches, such as Deep Reinforcement Learning [2], [3], [4], try to learn the control law using black-box methods, and hence, do not require any engineering for the specific system. These black-box methods abandon all insights from control and physics, require millions of samples from the physical systems, do not yield any performance guarantees

This project has received funding from the European Union's Horizon 2020 research and innovation program under grant agreement No #640554 (SKILLS4ROBOTS). Furthermore, this research was also supported by grants from ABB, NVIDIA and the NVIDIA DGX Station.

¹Michael Lutter and Jan Peters are with the Department of Computer Science, Technische Universität Darmstadt, 64289 Darmstadt, Germany {lutter, peters}@ias.tu-darmstadt.de

²Kim Listmann is with ABB Corporate Research Center Germany, Wallstadter Str. 59, 68526 Ladenburg, Germany kim.listmann@de.abb.com

³Jan Peters is with the Max Planck Institute for Intelligent Systems, Spemannstr. 41, 72076 Tübingen, Germany jan.peters@tuebingen.mpg.de



Fig. 1. The physical Cartpole and Furuta pendulum used to evaluate the energy-based control of under-actuated systems with learned models. (a) The Cartpole consists of a horizontally moving cart driven by a rack and pinion drive. The passive pendulum is attached to the cart and the cart the pendulum can be swung-up and balanced. The front (b) and top (c) view of the Furuta pendulum consisting of an actuated rotary pendulum with a passive vertical pendulum. The vertical pendulum can be swung-up and balanced by moving the horizontal rotary link. Videos for the swing-up using the different models are available at https://youtu.be/m3JRYq7Gmgo

and require extensive reward shaping to the desired solution [5] or random seeds [6].

We propose to bridge this gap by combining the flexibility of deep learning with the theoretical insights from control theory in order to achieve learning of control that is independent of the system, applicable for real systems, cannot yield degenerate solutions and requires little engineering. Therefore, we combine existing control laws for energy based control with model learning. Such combination cannot be achieved by standard black-box model learning techniques [7], [8], [9], [10] because these methods learn the mapping from joint state $\{\mathbf{q}, \dot{\mathbf{q}}, \ddot{\mathbf{q}}\}$ to motor torques $\boldsymbol{\tau}_{M}$, but cannot learn the underlying ODE¹ nor the potential and kinetic energies, because these components are not observable and hence, cannot be learned supervised. Only our novel extension of Deep Lagrangian Networks (DeLaN) [12] is capable of learning the underlying ODE from data using the joint configurations and motor torques. Compared to the previous DeLaN, we extend DeLaN to also encode energy conservation and coherence besides the Lagrangian Mechanics prior. Therefore, our novel extension of DeLaN learns the mass-matrix, the centrifugal, Coriolis, gravitational

¹One can also not infer the components of the ODE using the Composite Rigid body algorithm [11] in combination with the learned inverse dynamics mapping, because the black-box function approximation violates the underlying assumptions.

and frictional forces as well as the potential and kinetic energy using unsupervised learning. Hence, DeLaN enables the combination of energy control with learned models without the knowledge of the system kinematics, which must be known for standard system identification techniques [13]. In the following we will refer to this combination as DeLaN for energy control (DeLaN 4EC)

To demonstrate the performance, we apply DeLaN 4EC under-actuated systems. This control problem is challenging as the controller must exploit the inherent system dynamics to solve the task and cannot use high-gain feedback control to cancel the system dynamics. For example the swing-up task of the Cartpole and the Furuta pendulum requires the repeated amplification of the amplitude of the passive pendulum before the pendulum can be swung-up. Furthermore, these tasks are a standard evaluation task for learning for control [14]. In contrast to most previous research, we apply the control laws also to the physical Furuta pendulum (Figure 1) and learn the control without pre-training in simulation.

Contribution

The contribution of this paper is the novel extension of DeLaN and the combination of DeLaN and energy control (DeLaN 4EC) for controlling under-actuated systems. First, we extend DeLaN to incorporate energy conservation and frictional forces. Therefore, the extended DeLaN not only adheres to the Lagrangian Mechanics but also ensures energy conservation, temporal coherence of the energy and the passivity of the learned representation. Second, we demonstrate that this combination can achieve energy-control for the Cartpole and the Furuta pendulum. This is demonstrated in simulation and on the physical Furuta pendulum in real-time at 500Hz and without pre-training in simulation. The performance is compared to the analytic models of the manufacturer as well as the standard system identification approach [13].

In the following we provide an overview about related work (Section II), briefly summarize Deep Lagrangian Networks [12] and highlight the proposed extensions to this approach (Section III). Subsequently, we derive our proposed control approach, DeLaN for energy control (DeLaN 4EC) and state the energy-based control law (Section IV). Finally, the experiments in Section V evaluate the control performance for simulated and physical under-actuated systems.

II. RELATED WORK

Controlling under-actuated systems has been addressed from various perspectives including reinforcement learning and control theory. For reinforcement learning the swing-up of passive pendulums is a standard benchmark for continuous state- and action spaces. These methods learn the control policy by treating the control task as black-box and improve the policy using only scalar rewards as feedback. However, most reinforcement learning algorithms can only be used in simulation due to the high sample complexity. Only PILCO [15] learned the Cartpole swing-up on the physical system. From a control perspective many control laws for specific under-actuated systems have been proposed [16], [17], [18], [19]. These papers manually derive the dynamics for each system using Lagrangian Mechanics and use the specific equations to derive control laws. For the resulting control laws the stability can be analyzed and guaranteed given the true model [18]. Therefore, the control laws achieve the desired behaviour and cannot exploit ill-posed reward functions but require engineering of the dynamics and control law. With DeLaN 4EC we use the control perspective and embed a control law within a learning architecture to learn the complete control approach. Rather than using the specific system dynamics for deriving the control law, we use the generic Euler-Lagrange ODE, which describes any mechanical system including closed-loop kinematics, and learn the ODE describing the model from data.

Learning the model from data has been addressed in the literature by either system identification or supervised black-box function approximation. For system identification the knowledge of the kinematic structure is exploited such that the linkage physics parameters can be inferred using linear regression [13]. However, the learned parameters must not necessarily be physically plausible [20], can only be linear combinations and can only be applied to kinematic trees [21]. In combination with the composite rigid body algorithm [11] the parameters of the Euler-Lagrange ODE including the mass-matrix can be inferred. For the function approximations standard machine learning techniques such as Linear Regression [7], [22], Gaussian Mixture Regression [23], [24], Gaussian Process Regression [25], [9], [26], Support Vector Regression [8], [27], feedforward- [28], [29], [30], [10] or recurrent neural networks [31] have been used. These models learn the forward or inverse mapping from joint configuration $\{\mathbf{q}, \dot{\mathbf{q}}, \ddot{\mathbf{q}}\}$ to motor torque $\boldsymbol{\tau}_M$. Therefore, these learned models cannot be used to infer the parameters of the Euler-Lagrange ODE and do not allow a combination with classical control besides inverse dynamics or non-linear feed-forward control [26].

In contrast to these existing methods, DeLaN learns the Euler-Lagrange ODE directly from data, does not require any knowledge of the kinematic structure and is not restricted to kinematic trees. Therefore, DeLaN learns the mass matrix, the centrifugal-, Coriolis-, gravitational- and frictional forces as well as the system energy using unsupervised learning and fits naturally with control theory.

III. DEEP LAGRANGIAN NETWORKS

First in Section III-A, the concept of Deep Lagrangian Networks [12] is summarized and novel extensions are proposed in the subsequent sections. Section III-B extends the cost function with the forward model. Section III-C introduces friction such that the Lagrangian Mechanics prior is not violated. Finally, section III-D adds energy conservation as additional constraint to model learning. Thus, the extended DeLaN not only complies with Lagrangian Mechanics but also ensures energy conservation and coherence.



Fig. 2. The computational graph of the Deep Lagrangian Network for control (DeLaN 4EC). Shown in blue and green is the neural network with the three separate heads computing the potential energy V and the mass-matrix **H**. The orange boxes construct represent the physics transformations constructing Euler-Lagrange equation. For energy-based control these components are directly interfaced to the control law to determine the motor-torque. For training, the gradients are backpropagated through all vertices highlighted in orange.

A. Deep Lagrangian Networks

Deep Lagrangian Networks use the knowledge from Lagrangian Mechanics and encode this prior within a deep learning architecture. Therefore, all learned models guarantee that these models must comply with Lagrangian Mechanics. More concretely, let the Lagrangian be defined as L = T - V, where $T = 1/2 \dot{\mathbf{q}}^T \mathbf{H}(\mathbf{q}) \dot{\mathbf{q}}$ is the kinetic energy, V the potential energy and **H** the positive definite mass matrix. Substituting L into the Euler-Lagrange differential equation yields the ODE described by

$$\mathbf{H}(\mathbf{q})\ddot{\mathbf{q}} + \underbrace{\dot{\mathbf{H}}(\mathbf{q})\dot{\mathbf{q}} - \frac{1}{2} \left(\frac{\partial}{\partial \mathbf{q}} \left(\dot{\mathbf{q}}^T \mathbf{H}(\mathbf{q})\dot{\mathbf{q}} \right) \right)^T}_{:=\mathbf{C}(\mathbf{q},\dot{\mathbf{q}})\dot{\mathbf{q}}} + \frac{\partial V}{\partial \mathbf{q}} = \sum_i \boldsymbol{\tau}_i \quad (1)$$

where $\boldsymbol{\tau}_i$ are the non-conservative generalized forces including motor and frictional forces. Approximating **H** and *V* using deep networks, i.e.,

$$\hat{\mathbf{H}} = \hat{\mathbf{L}}(\mathbf{q};\theta) \, \hat{\mathbf{L}}^{T}(\mathbf{q};\theta) + \epsilon \, \mathbf{I} \qquad \hat{V} = \hat{V}(\mathbf{q};\psi) \tag{2}$$

where $\hat{\cdot}$ refers to an approximation, $\hat{\mathbf{L}}$ is a lower triangular matrix with a non-negative diagonal, θ and ψ are the network parameters and ϵ a small positive constant, one can encode the ODE by exploiting the full differentiability of the neural networks [12]. Additionally, the mass matrix **H** is guaranteed to be positive definite and the eigenvalues are lower-bounded by ϵ . The network parameters can be learned online and end-to-end, by minimizing the error of the ODE using the samples { $\mathbf{q}, \dot{\mathbf{q}}, \mathbf{\tau}_M$ } recorded on the physical system, i.e. minimizing the ℓ_i norm between the prediction of Equation 1 and the observed motor torque $\boldsymbol{\tau}_M$. Therefore, the superposition of the different forces is learned supervised, while the decomposition into inertial, Coriolis, centripetal and gravitational forces is learned unsupervised.

B. Introducing the Forward Model

Unlike many model learning techniques, DeLaN can be used as forward and inverse dynamics model, by solving Equation 1 w.r.t. $\ddot{\mathbf{q}}$. Therefore, one can incorporate the loss of the forward model within the learning of the parameters. This is especially important for many control approaches, including energy control, as these use the inverse of the mass matrix. Therefore, incorporating the forward model within the learning of the parameters should yield better approximation of the inverse. Solving Equation 1 for $\ddot{\mathbf{q}}$ yields

$$\mathbf{H}^{-1}(\mathbf{q})\left(\sum_{i}\boldsymbol{\tau}_{i}-\dot{\mathbf{H}}(\mathbf{q})\dot{\mathbf{q}}+\frac{1}{2}\left(\frac{\partial}{\partial \mathbf{q}}\left(\dot{\mathbf{q}}^{T}\mathbf{H}(\mathbf{q})\dot{\mathbf{q}}\right)\right)^{T}-\frac{\partial V}{\partial \mathbf{q}}\right)=\ddot{\mathbf{q}}.$$

Thus the loss function can be extended to minimize the error of the inverse and forward model, i.e.,

$$(\theta^*, \psi^*) = \underset{\theta, \psi}{\operatorname{arg\,min}} \ell_i \left(\hat{f}(\theta, \psi), \ \ddot{\mathbf{q}} \right) + \\ \ell_i \left(\hat{f}^{-1}(\theta, \psi), \ \boldsymbol{\tau}_M \right) + \lambda \Omega(\theta, \psi)$$

$$(3)$$

where Ω is the l_2 weight regularization.

C. Introducing Friction to Model Learning

Incorporating friction within model learning in a non black-box fashion is non-trivial because friction is an abstraction to combine various physical effects. For robot arms in free space the friction of the motors dominates, for mechanical systems dragging along a surface the friction at surface dominates while for legged locomotion the friction between the feet and floor dominates but also varies with time. Therefore, defining a general case for all types of friction in compliance with the Lagrangian Mechanics is challenging. Various approaches to incorporate friction models can be found in [32], [33]. Furthermore, if the friction model includes stiction the dynamics are not invertible because multiple motor-torques can generate the same joint acceleration [34].

This paper focuses on friction caused by the actuators. For actuator friction different models have been proposed [35], [1], [36], [37]. These models assume that the motor friction only depends on the joint velocity $\dot{\mathbf{q}}_i$ of the *i*th-joint and is independent of the other joints [35], [1], [36], [37]. Depending on model complexity a combination of static, viscous or Stribeck friction is assumed as model prior and the superposition is described by

$$\boldsymbol{\tau}_{f_i} = -\left(\tau_{C_v} + \tau_{C_s} \exp\left(-\dot{\mathbf{q}}_i^2/\nu\right)\right) \operatorname{sign}(\dot{\mathbf{q}}_i) - d\,\dot{\mathbf{q}}_i \qquad (4)$$

where τ_{C_v} is the coefficient of static friction, *d* the coefficient of viscous friction, and τ_{C_s} and ν are the coefficients of Stribeck friction. In the following the friction coefficients are abbreviated as $\phi = \{\tau_{C_v}, \tau_{C_s}, \nu, d\}$. Since the frictional force τ_f is a function of the generalized coordinates, the frictional force is a non-conservative and generalized force and can simply be added to the Lagrange Euler ODE (Equation 1). For other types of friction this is not true and one needs to explicitly ensure that one can express the frictional force as generalized force. Given the model prior of Equation 4 the friction coefficients ϕ can be learned by treating the coefficients as network weights.

D. Introducing Energy to Model Learning

Besides the Lagrangian Mechanics objective, incorporating energy conservation and energy coherence, i.e., ensuring that $E_i(t) \forall t \ge 0$ is at least of Class C^2 , is natural because the Lagrangian L contains the system energy. In order to ensure the conservation of energy, the total energy of the system must be equal to the summation of the initial system energy E_0 , the work done by the actuators W_m and the energy losses due to friction E_{th} , i.e.,

$$E(t) = T(t) + V(t) = E_0 + W_M(t) + E_{th}(t) \quad \forall t \ge 0.$$
 (5)

The actuator work or the losses to friction can be computed by numerical integration described by

$$W_j(t) = \int_{\mathbf{q}(0)}^{\mathbf{q}(t)} \boldsymbol{\tau}_j^T(\mathbf{q}) \ d\mathbf{q} = \int_0^t \boldsymbol{\tau}_j^T(\mathbf{q}(u)) \ \dot{\mathbf{q}}(u) \ du \qquad (6)$$

where τ_j is either the frictional torque τ_F or actuator torque τ_M and $W_j(0) \coloneqq 0$. This can also be expressed in using the change in energy, i.e.,

$$\dot{E} = \dot{\mathbf{q}}^{T} \left(\boldsymbol{\tau}_{M} + \boldsymbol{\tau}_{F} \right) = \dot{T} + \dot{V} = \dot{\mathbf{q}}^{T} \mathbf{H} \ddot{\mathbf{q}} + \frac{1}{2} \dot{\mathbf{q}}^{T} \dot{\mathbf{H}} \dot{\mathbf{q}} + \dot{\mathbf{q}}^{T} \frac{\partial V}{\partial \mathbf{q}}.$$
(7)

Following [21] and recognizing that $\tau_M + \tau_F$ is the total force acting on the mechanical system, Equation 7 not only ensures energy conservation but also the passivity of the learned system, because this equality ensures the lower bound on the total energy, i.e., $E(T) - E(0) \ge -E(0) \forall T > 0$. Therefore, the learned model representation is guaranteed to be passive on the training domain given sufficiently low training error. This property of DeLaN implies that the uncontrolled system described by the learned dynamics is stable. For black-box function approximation methods this must not be necessarily be true because these methods can learn an active system that is optimal w.r.t. the given cost function. Besides the conservation of energy, the energy coherence can be used as additional constraint ensuring that both the kinetic and potential energy is continuous and differentiable w.r.t. time, i.e., $T, V \in C^1$. Using a first order Taylor approximation this constraint can be expressed as

$$\tilde{T}(\mathbf{q}_{t+\delta t};\theta) = \hat{T}(\mathbf{q}_t;\psi) + \dot{\mathbf{q}}_t^T \mathbf{H} \ddot{\mathbf{q}}_t \delta_t + \frac{1}{2} \dot{\mathbf{q}}_t^T \dot{\mathbf{H}} \dot{\mathbf{q}}_t^T \delta_t \quad (8)$$

$$\tilde{V}(\mathbf{q}_{t+\delta t};\psi) = \hat{V}(\mathbf{q}_t;\psi) + \dot{\mathbf{q}}_t^T \frac{\partial V}{\partial \mathbf{q}} \delta_t.$$
(9)

The resulting equations cannot be directly used as a loss because the true kinetic- and potential energy of the configuration $\mathbf{q}_{t+\delta_t}$ is unknown. Therefore, we bootstrap the current approximation of \tilde{V} and \tilde{T} as target value and do not propagate the gradients through these estimates. In addition, the energy for a specific joint configuration \mathbf{q}^* is clamped to a pre-specified value as in [38], i.e. $E(\mathbf{q}^*, \dot{\mathbf{q}}^*) \coloneqq 0$. Adding energy conservation (Equation 7) and energy coherence (Equation 8 & Equation 9) to the optimization problem of Equation 3 yields the loss for DeLaN 4EC.

IV. DEEP LAGRANGIAN NETWORKS FOR ENERGY CONTROL

In the previous section, we showed that DeLaN can learn the mass matrix, the centripetal, gravitational and frictional forces as well as the kinetic and potential energies using only the joint measurements $(\mathbf{q}, \dot{\mathbf{q}}, \ddot{\mathbf{q}})$ and the actuator torques $\boldsymbol{\tau}_{M}$. Using these properties, energy control can be achieved by embedding the learned energies within a energy-based control law. Therefore, DeLaN 4EC enables the control of a large-class of under-actuated systems, because these systems are mainly controlled using energy-based control laws and other black-box identification techniques cannot learn the system energy and hence, cannot be applied. For energy control, the control law proposed by Spong et. al. [17], which is applicable to the Furuta pendulum, the Cartpole and the Acrobot is used. This control law regulates the energy of the pendulum E_p to obtain the desired energy E^* and adds an additional P-controller on the active joints to avoid the joint limits. For systems with high friction an additional term to compensate the friction of the actuator can be added. Expressing this control law using the mass-matrix and the potential energy is described by

$$\mathbf{u}_{E} = k_{E} \left(E_{p} - E^{*} \right) \operatorname{sign} \left(\dot{\mathbf{q}}_{p} \cos(\mathbf{q}_{p}) + \mathbf{K}_{p} \left(\mathbf{q}_{a}^{*} - \mathbf{q}_{a} \right) \quad (10)$$

with the pendulum energy $E_p = 1/8 \dot{\mathbf{q}}_p^T \mathbf{H}_{22} \dot{\mathbf{q}}_p + V(\mathbf{q})$ and the desired energy $E^*(\mathbf{q}^*, \dot{\mathbf{q}}^*)$ at the desired joint configuration \mathbf{q}^* .

V. EXPERIMENTS

We apply DeLaN 4EC to control two under-actuated systems: the Cartpole (Figure 1a) and the Furuta pendulum (Figure 1b). The Cartpole is a horizontally moving cart with an attached passive pendulum. Moving the cart horizontally indirectly controls the pendulum and using this indirect control the pendulum can be swung-up and balanced. Similarly, the Furuta pendulum (also referred to as whirling pendulum) consists of an actuated rotary pendulum with a vertical passive pendulum. Using the rotary pendulum the vertical link can be swung-up and balanced. These experiments are standard experiments for learning to control. However, most previous research only used simulations while we apply these methods to the physical Cartpole and Furuta pendulum.

A. Experimental Setup

To learn the control task, a smooth exploration policy, i.e., the energy controller using the analytic model, interacts for T seconds with the system and generates data containing



Fig. 3. The motor-, centrifugal-, Coriolis-, gravitational- and frictional forces, the joint acceleration as well as the kinetic and potential energy for the swing-up of the simulated and physical Cartpole. Using only the super-imposed motor torques and joint accelerations as supervising feedback, DeLaN learns to disambiguate between the individual force components and system energies.



Fig. 4. The normalized mean squared error of the forward and inverse models for the simulated and physical platforms on the test set.

{q, \dot{q} , \ddot{q} , τ_M }_{1...N}. For the Furuta pendulum the interaction time is 120s while the interaction time for the Cartpole is 240s. Using these highly correlated samples the control is learned offline. After learning the controller, the performance is evaluated using the normalized mean square error (nMSE) on the test data (Section V-B) as well as the online control performance on the swing-up tasks (Section V-C). The online control evaluation is the more relevant performance measure as the nMSE can be deceiving and a low nMSE does not necessarily imply a good control performance. For the swing-

up task the systems are first stabilized to the desired energy E^* of the balancing point using energy control and then balanced at the unstable equilibrium using a PD-controller. Both, controller operate at 500Hz and the gains are tuned for each system using the analytic model provided by the manufacturer and fixed afterwards for each experiment.

The simulated experiments are performed using Bullet [39] with joint torque as control input. The physical experiments are performed using the Cartpole (Fig. 1a) and Furuta pendulum (Fig. 1b) manufactured by Quanser. These physical systems are directly controlled using the DC motor voltage. For the experiments the voltage to motor-torque conversion is performed using the parameters of the manufacturer. Furthermore, both physical systems have unique properties that make the model learning and the control challenging. The linear actuation of the Cartpole is a pinion & rack drive causing significant stiction and this stiction renders the model learning challenging. In contrast, the links of the Furuta pendulum are very light weight and even small errors of motor voltages push the active joints to its joint limit and stop the episode.

The performance of DeLaN 4EC is compared to the dynamics parameters of the manufacturer and the white-box system identification introduced by [13] with the extension of viscous friction as in [20]. For system identification the mass matrix is computed using the Composite Rigid Body algorithm [11] and the potential energy is computed using the analytic expression $V(\mathbf{q}) = mgl(\cos(\mathbf{q}) + 1)$, where only the mass *m* is inferred from data while the gravitational constant *g* and pendulum length *l* are pre-defined con-



Fig. 5. The position θ and velocity $\dot{\theta}$ orbits recorded using energy control to swing-up the passive pendulums. The rows show the different models, i.e., the analytic model, the system identification model and the DeLaN model while the columns show the different simulated and physical systems. The dahsed orbit highlights the desired energy E^* . While the learned and the analytic model can swing-up the simulated system and physical Cartpole only the analytic model and DeLaN 4EC can swing-up the physical Furuta pendulum, while the energy controller using the System Identification model cannot.

stants. These requirements are in stark contrast to DeLaN 4EC, because the white-box system identification approach requires the kinematic structure defining the link length, connection between links and gravitational constant, while DeLaN 4EC must learn the kinematic and dynamic structure from data. Furthermore, the assumption of knowing the kinematic structure simplifies the learning of the potential energy to merely fitting the amplitude of the potential energy while DeLaN 4EC must not only learn the amplitude but also learn the shape. We do not compare to other blackbox learning techniques such as neural networks or Gaussian process regression because these techniques cannot learn energy and hence, cannot be applied to energy control.

B. Offline Evaluation

Figure 3 shows the learned dynamics model and energies of the simulated and physical Cartpole executing a swing-up using the control law described by \mathbf{u}_E . The dynamics learned by DeLaN closely resemble the data as well as the frictional-, inertial-, centrifugal-, Coriolis- and gravitational forces predicted by the analytic model. Furthermore, the kinetic and potential energies are learned. Only the learned potential and kinetic energy for the physical system are slightly scaled similar to the system identification model. Although the energies are scaled, both the kinetic and potential energy are scaled coherently such that the energy conservation holds as enforced by Equation 9. Furthermore, DeLaN learns the high stiction $\boldsymbol{\tau}_F$ of the pinion and rack drive and predicts the close to zero accelerations $\ddot{\mathbf{q}}$ and non-zero motor torques $\boldsymbol{\tau}_{M}$ during balancing of the physical Cartpole, given a sufficiently good initialization of the friction model. The analytic model and system identification cannot represent the stiction and predict either zero torques and non-zero accelerations or zero torques and zero acceleration. Both predictions oppose the measured data. However, DeLaN suffers from high frequency noise on the passive pendulum or close to zero components, whereas the white box models do not because the white box models consist of global parameters, which are not susceptible to noise and these models can exploit the known kinematic structure to infer zero Coriolis or gravitational forces.

Figure 4 shows the quantitative comparison using the nMSE defined as

$$nMSE = \frac{\sum_{i=0}^{N} \|\mathbf{x}_{i} - \hat{\mathbf{x}}_{i}\|_{2}^{2}}{\sum_{i=0}^{N} \|\mathbf{x}_{i} + \delta\|_{2}^{2}}$$
(11)

whereas δ is a small constant for numerical stability. The nMSE is evaluated on test data performing a swing-up, which in the case of the Cartpole is identical to the data shown in Figure 3. For the simulations the analytic model is the true

model, which has a non-zero nMSE because noise is added to the torques during simulation and the accelerations are computed using finite differences and are low-pass filtered because this signal-processing is required for the physical system. The comparison shows that DeLaN obtains a similar nMSE as system identification and the analytic model for the forward model of the simulated systems. For the inverse model of the simulated systems, DeLaN obtains comparable nMSE for the actuated joint but slightly increased nMSE for the passive joints because DeLaN is susceptible to noise and the nMSE is very sensitive to the noise of the passive joint as $\tau_p \coloneqq 0$. For the physical systems, DeLaN and system identification obtain a lower nMSE than the analytic model for the forward model. For the inverse model, both learned models achieve better performance than the analytic model on the active joint and only for the passive joint DeLaN performs slightly worse due to the noise. This noise is negligible during optimization and control because the MSE is dominated by the actuated joint and only the nMSE per joint amplifies the impact of the noise. Overall, the qualitative and quantitative evaluation showed that the performance of the learned models in comparison to the analytic model achieve comparable performance for the simulated systems and a slightly better performance for the physical systems.

TABLE I Percentage of successful swing-ups of simulated and physical Cartpole and Furuta pendulum for the different models.

	Car	Cartpole		Furuta	
Model	Sim	Real	Sim	Real	
Analytic Model	1.00	1.00	1.00	1.00	
System Identification	1.00	1.00	0.93	0.00	
DeLaN	1.00	1.00	1.00	0.90	

C. Online Control Evaluation

For the online control experiments the energy-based control law described in Equation 10 is applied with a control frequency of 500Hz to 30 different initial joint configurations. Therefore, all models must achieve real-time computation of at least 500Hz on the physical system to be able to solve the task. For the simulated experiments, the starting configuration is randomly sampled while the physical experiments are performed sequentially and hence, the starting configuration naturally changes. For the physical Cartpole we augment the energy controller with an negative derivative gain to compensate the large viscous friction of the pinion and rack drive. The percentage of successful swingups is summarized in Table I and the corresponding position θ and velocity $\dot{\theta}$ orbits for two starting configurations of the passive pendulum are shown in Figure 5. Videos of the swing-up of the physical Cartpole and physical Furuta pendulum can be found at https://youtu.be/m3JRYq7Gmgo.

For the simulated systems the analytic and the learned models achieve the successful completion of the swingup. Only the system identification model fails on 2 trials. These unsuccessful completions are caused by the balancing controller because the system identification model swingsup the pendulum with slightly too much or too low energy such that the PD-controller fails to stabilize the pendulum. Furthermore, the resulting trajectories for the learned models are indistinguishable. For the physical Cartpole all models achieve smooth real-time control and swing-up the pendulum for 30 consecutive times from varying starting configurations. For the physical Furuta pendulum only the analytic model and DeLaN 4EC achieve the successful swing-up, while the system identification model can only stabilize the pendulum to a low amplitude cycle, which does not reach the balancing point. For 30 trials, DeLaN 4EC achieves the successful completion for 27 trials. The three unsuccessful trials are caused by the PD-Controller not being able to stabilize the pendulum because DeLaN 4EC swings up the pendulum with slightly less energy than the analytic model and the balancing PD-Controller is very sensitive to these changes in velocity at the switching point. Tuning the PD-controller gains w.r.t. to DeLaN 4EC results in the successful completion of all 30 trials but decreases the performance of the analytic model. For fair comparison the gains were fixed between the experiments and optimized for the analytic model. The system identification model fails to swing-up the pendulum because this approach learns a too low mass for the pendulum and hence, can only stabilize the pendulum to a low amplitude oscillation. The learning fails because the regressor of the system identification has too low rank and can only infer a linear combination of the dynamics parameters [40].

DeLaN 4EC is capable of solving the swing-up for the simulated and physical Cartpole and Furuta pendulum using a 500Hz real-time control loop. The performance of DeLaN 4EC is comparable to the analytic model and DeLaN 4EC achieves the swing-up of the physical Furuta pendulum, where system identification does not, despite having a lower nMSE compared to the analytic model. This shows that a low nMSE does not necessarily imply good control performance.

VI. CONCLUSION

In this paper, we introduced the concept of Deep Lagrangian Networks for energy control (DeLaN 4EC), a learning to control approach that combines the flexibility of deep learning with the insights from control theory. This combination is enabled only because DeLaN 4EC imposes Lagrangian Mechanics, conservation of energy and energy coherence on a generic deep network and hence, learns a physically plausible model. We showed that DeLaN is able to learn the inertial-, centripetal-, Coriolis-, gravitationaland frictional forces and the potential and kinetic energy from sensor data containing only joint configuration and motor torque. Therefore, learning these forces and energies is unsupervised and does not require any knowledge about the kinematic structure. Other model learning algorithms either require the kinematic structure to learn these components such as system identification or cannot learn the force components or system energies such as neural networks or Gaussian process regression. The qualitative and quantitative offline evaluation showed that the normalized MSE of DeLaN in comparison to the analytic model is comparable for the simulated systems and better for the physical systems. For the online control task, DeLaN 4EC accomplishes the swing-up of the physical Cartpole and Furuta pendulum from different starting configurations by computing the system energies within a 500Hz real-time control loop. In contrast, the system identification model only achieves the successful swing-up of the physical Cartpole but not the physical Furuta pendulum, despite having comparable nMSE to DeLaN 4EC.

References

- A. Albu-Schäffer, "Regelung von Robotern mit elastischen Gelenken am Beispiel der DLR-Leichtbauarme," Ph.D. dissertation, Technische Universität München, 2002.
- [2] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," arXiv preprint arXiv:1509.02971, 2015.
- [3] J. Schulman, S. Levine, P. Abbeel, M. I. Jordan, and P. Moritz, "Trust region policy optimization." in *Icml*, vol. 37, 2015, pp. 1889–1897.
- [4] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," arXiv preprint arXiv:1707.06347, 2017.
- [5] M. J. Mataric, "Reward functions for accelerated learning," in *Machine Learning Proceedings 1994*. Elsevier, 1994, pp. 181–189.
- [6] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger, "Deep reinforcement learning that matters," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [7] S. Schaal, C. G. Atkeson, and S. Vijayakumar, "Scalable techniques from nonparametric statistics for real time robot learning," *Applied Intelligence*, vol. 17, no. 1, pp. 49–60, 2002.
- [8] Y. Choi, S.-Y. Cheong, and N. Schweighofer, "Local online support vector regression for learning control," in *International Symposium on Computational Intelligence in Robotics and Automation*. IEEE, 2007, pp. 13–18.
- [9] D. Nguyen-Tuong, M. Seeger, and J. Peters, "Model learning with local gaussian process regression," *Advanced Robotics*, vol. 23, no. 15, pp. 2015–2034, 2009.
- [10] A. Sanchez-Gonzalez, N. Heess, J. T. Springenberg, J. Merel, M. Riedmiller, R. Hadsell, and P. Battaglia, "Graph networks as learnable physics engines for inference and control," *arXiv preprint arXiv:1806.01242*, 2018.
- [11] M. W. Walker and D. E. Orin, "Efficient dynamic computer simulation of robotic mechanisms," *Journal of Dynamic Systems, Measurement, and Control*, vol. 104, no. 3, pp. 205–211, 1982.
- [12] M. Lutter, C. Ritter, and J. Peters, "Deep lagrangian networks: Using physics as model prior for deep learning," in *International Conference* on Learning Representations, 2019.
- [13] C. G. Atkeson, C. H. An, and J. M. Hollerbach, "Estimation of inertial parameters of manipulator loads and links," *The International Journal* of Robotics Research, vol. 5, no. 3, pp. 101–119, 1986.
- [14] Y. Duan, X. Chen, R. Houthooft, J. Schulman, and P. Abbeel, "Benchmarking deep reinforcement learning for continuous control," in *International Conference on Machine Learning*, 2016, pp. 1329– 1338.
- [15] M. Deisenroth and C. E. Rasmussen, "Pilco: A model-based and data-efficient approach to policy search," in *Proceedings of the 28th International Conference on machine learning (ICML-11)*, 2011, pp. 465–472.
- [16] C. C. Chung and J. Hauser, "Nonlinear control of a swinging pendulum," *automatica*, vol. 31, no. 6, pp. 851–862, 1995.
- [17] M. W. Spong, "Energy based control of a class of underactuated mechanical systems," *IFAC Proceedings Volumes*, vol. 29, no. 1, pp. 2828–2832, 1996.
- [18] I. Fantoni and R. Lozano, Non-linear control for underactuated mechanical systems. Springer Science & Business Media, 2002.
- [19] M. Ishitobi, Y. Ohta, Y. Nishioka, and H. Kinoshita, "Swing-up of a cart—pendulum system with friction by energy control," *Proceedings* of the Institution of Mechanical Engineers, Part I: Journal of Systems and Control Engineering, vol. 218, no. 5, pp. 411–415, 2004.
- [20] J.-A. Ting, M. Mistry, J. Peters, S. Schaal, and J. Nakanishi, "A bayesian approach to nonlinear parameter identification for rigid body dynamics." in *Robotics: Science and Systems*, 2006, pp. 32–39.

- [21] M. W. Spong, S. Hutchinson, M. Vidyasagar et al., Robot modeling and control. Wiley New York, 2006, vol. 3.
- [22] M. Haruno, D. M. Wolpert, and M. Kawato, "Mosaic model for sensorimotor learning and control," *Neural computation*, vol. 13, no. 10, pp. 2201–2220, 2001.
- [23] S. Calinon, F. D'halluin, E. L. Sauser, D. G. Caldwell, and A. G. Billard, "Learning and reproduction of gestures by imitation," *IEEE Robotics & Automation Magazine*, vol. 17, no. 2, pp. 44–54, 2010.
- [24] S. M. Khansari-Zadeh and A. Billard, "Learning stable nonlinear dynamical systems with gaussian mixture models," *IEEE Transactions* on *Robotics*, vol. 27, no. 5, pp. 943–957, 2011.
- [25] J. Kocijan, R. Murray-Smith, C. E. Rasmussen, and A. Girard, "Gaussian process model based predictive control," in *American Control Conference*, vol. 3. IEEE, 2004, pp. 2214–2219.
- [26] D. Nguyen-Tuong and J. Peters, "Using model knowledge for learning inverse dynamics." in *International Conference on Robotics and Automation*, 2010, pp. 2677–2682.
- [27] J. P. Ferreira, M. Crisostomo, A. P. Coimbra, and B. Ribeiro, "Simulation control of a biped robot with support vector regression," in *IEEE International Symposium on Intelligent Signal Processing*. IEEE, 2007, pp. 1–6.
- [28] M. Jansen, "Learning an accurate neural model of the dynamics of a typical industrial robot," in *International Conference on Artificial Neural Networks*, 1994, pp. 1257–1260.
- [29] I. Lenz, R. A. Knepper, and A. Saxena, "Deepmpc: Learning deep latent features for model predictive control." in *Robotics: Science and Systems*, 2015.
- [30] F. D. Ledezma and S. Haddadin, "First-order-principles-based constructive network topologies: An application to robot inverse dynamics," in *IEEE-RAS International Conference on Humanoid Robotics*, 2017. IEEE, 2017, pp. 438–445.
- [31] E. Rueckert, M. Nakatenus, S. Tosatto, and J. Peters, "Learning inverse dynamics models in o (n) time with lstm networks," in *IEEE-RAS International Conference on Humanoid Robotics*. IEEE, 2017, pp. 811–816.
- [32] A. I. Lurie, Analytical mechanics. Berlin, Heidelberg: Springer Science & Business Media, 2013.
- [33] D. A. Wells, Schaum's outline of theory and problems of lagrangian dynamics. McGraw-Hill, 1967.
- [34] N. Ratliff, F. Meier, D. Kappler, and S. Schaal, "Doomed: Direct online optimization of modeling errors in dynamics," *Big data*, vol. 4, no. 4, pp. 253–268, 2016.
- [35] H. Olsson, K. J. Åström, C. C. De Wit, M. Gäfvert, and P. Lischinsky, "Friction models and friction compensation," *Eur. J. Control*, vol. 4, no. 3, pp. 176–195, 1998.
- [36] B. Bona and M. Indri, "Friction compensation in robotics: an overview," in *Decision and Control*, 2005 and 2005 European Control Conference. CDC-ECC'05. 44th IEEE Conference on. IEEE, 2005, pp. 4360–4367.
- [37] A. Wahrburg, J. Bös, K. D. Listmann, F. Dai, B. Matthias, and H. Ding, "Motor-current-based estimation of cartesian contact forces and torques for robotic manipulators and its application to force control," *IEEE Transactions on Automation Science and Engineering*, vol. 15, no. 2, pp. 879–886, 2018.
- [38] M. Riedmiller, "Neural fitted q iteration-first experiences with a data efficient neural reinforcement learning method," in *European Conference on Machine Learning*. Springer, 2005, pp. 317–328.
- [39] E. Coumans and Y. Bai, "Pybullet, a python module for physics simulation for games, robotics and machine learning," http://pybullet.org, 2016–2018.
- [40] B. Siciliano and O. Khatib, Springer handbook of robotics. Springer, 2016.