Advances in Quantum Reinforcement Learning

Vedran Dunjko*
Institute for Theoretical Physics,
University of Innsbruck,
Innsbruck 6020, Austria
now at:

Max Planck Institute of Quantum Optics Garching 85748, Germany Email: vedran.dunjko@mpq.mpg.de Joint Quantum Institute & NIST, College Park, Maryland 20742, USA Email: jmtaylor@umd.edu

Jacob M. Taylor

Hans J. Briegel
Institute for Theoretical Physics,
University of Innsbruck
Innsbruck 6020, Austria, and
Department of Philosophy,
University of Konstanz,
Konstanz 78457, Germany
Email: hans.briegel@uibk.ac.at

Abstract-In recent times, there has been much interest in quantum enhancements of machine learning, specifically in the context of data mining and analysis. Reinforcement learning, an interactive form of learning, is, in turn, vital in artificial intelligence-type applications. Also in this case, quantum mechanics was shown to be useful, in certain instances. Here, we elucidate these results, and show that quantum enhancements can be achieved in a new setting: the setting of learning models which learn how to improve themselves - that is, those that metalearn. While not all learning models meta-learn, all non-trivial models have the potential of being "lifted", enhanced, to metalearning models. Our results show that also such models can be quantum-enhanced to make even better learners. In parallel, we address one of the bottlenecks of current quantum reinforcement learning approaches: the need for so-called oracularized variants of task environments. Here we elaborate on a method which realizes these variants, with minimal changes in the setting, and with no corruption of the operative specification of the environments. This result may be important in near-term experimental demonstrations of quantum reinforcement learning.

I. INTRODUCTION

¹The interest in the potential of quantum enhancements in aspects of machine learning (ML) and artificial intelligence (AI) has been on the rise in recent times. For example, quantum algorithms which can be used for efficient classification and clustering [1], [2], [3], [4], have been proposed. Such classification and clustering tasks, more generally referred to as supervised and unsupervised learning, form two out of three canonical branches of ML. These learning modes are complemented by reinforcement learning (RL) [5]. RL is central to some of the most celebrated recent successes in ML (e.g. a computer beating a human champion in the game of Go [6]), and is indispensable in the context of AI [7]. Nonetheless, it has received significantly less attention from the quantum information processing (QIP) community, in

¹©2017 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works. Full citation:

V. Dunjko, J. M. Taylor, H. J. Briegel, Advances in quantum reinforcement learning, IEEE SMC, Banff, AB, 2017, pp. 282-287. doi: 10.1109/SMC.2017.8122616 (2017).

comparison to the (un)supervised learning settings. Few works do exist: for instance in [8] it was suggested QIP may yield new ways to evaluate policies, and in [9] it was proven that the computational complexity of a particular model, Projective Simulation [10], can be quadratically reduced. Arguably, one of the key reasons why RL has been less explored stems from important differences to (un)supervised learning, which significantly influence how, and to what extent, QIP techniques can be applied. In particular, standard (un)supervised learning methods deal with a stationary dataset – a sample set from an (empirical) joint distribution P(X=x,Y=y), giving data points $(\{x_i\})$ and corresponding classes $(\{y_i\})$ in the supervised setting, or just the bare points sampled from P(X=x)in the unsupervised case. For convenience, we will abbreviate (un)supervised learning with DL (for data-learning). In contrast, RL is about interactive, reward-driven learning, involving a task environment. The most prominent feature of RL is that the learner, typically called the *learning agent*, influences the state of the environment - thereby changing the distribution of what it perceives. This distinction makes the RL model particularly well-suited for AI-type applications (rather than data analysis). The same distinction causes additional technical obstacles in defining quantum analogs, or enhancements, for RL, relative to DL. To clarify, since in DL we are dealing with a stationary database, we can think of it as being encoded in a memory where data points are accessible via an address (such as RAM or hard and flash drives). It will be convenient to think of such memory as a process, which, given an address or index, outputs the corresponding data point: $i \stackrel{Mem}{\rightarrow} x_i$. Such a memory can easily be represented as a reversible process, and consequently as a quantum operation, a unitary map, performing $|i\rangle |0\rangle \stackrel{qMem}{\rightarrow} |i\rangle |x_i\rangle$. In [11] it was shown how such a unitary can be implemented efficiently. Having access to such an object often allows for the efficient encoding of data into amplitudes, which is at the basis of many exponential quantum improvements [4], [21], or directly with the purpose of quadratic improvements [3]. The memory can be thought of as a black-box oracle and quantum-parallel access allows for query-complexity improvements, a topic which has been exhaustively explored from the early days of quantum computation. Note, in such settings, overall improvements are only attainable assuming such a memory has already been "loaded" with the data-set off-line, barring few exceptions [22]. In the RL case, however, the environment cannot be modelled as a stationary memory, as the environmental responses depend on the interaction history. In particular, they depend on the actions of the agent. The environment is thus a process with a memory of its own, which cannot be accessed by the agent – a memory with memory, so to speak. The general framework for quantum (RL) agents which interact with quantum environments was provided first in [12] and in the technical preprint [13] parts of which we summarize and expand upon here.

In the most general setting, both the learner and the environment are maps with memory (often called quantum combs [14]), and classical RL corresponds to the case where the maps are themselves classical. Any improvement beyond that of computational complexity requires that the environment allows some kind of quantum coherent access. This is an intuitive observation (and formally proven in [13]), in the same sense a classical phone book cannot be "Groverized" – it has to be first be mapped to a system which can act on superpositions. Thus this claim also applies for the case of DL, where the suitable quantum analog is straight-forward – it is the map qMem, see Fig 1. For the RL case, the situation is much less clear. In [12] it was shown that a quantum analog, which allows for provable learning improvements, does exist. In this work, we will expand on these results two-fold. We will show how such a map can be black-box constructed given any specification of a task environment, and consider cases when such a construction is possible. Furthermore, we will identify a new family of improvements, based on so-called meta-learning, which are possible given such an oracular construction. To do so, we first present the broad framework, outlined in [12], and the ideas of the key results we will need later.

A. Classical and quantum agent-environment interaction

The standard turn-based RL paradigm comprises the percept (\mathcal{S}) and action (\mathcal{A}) sets which specify the possible outputs of the environment, and the agent, respectively. The agent and the environment interact by sequentially exchanging elements from the percept/action sets. A realized interaction up to time step t, between the agent and the environment, that is a sequence $h_t = (s_1, a_2, s_3, s_4, \ldots, s_{t-1}, a_t), s_i \in \mathcal{S}, a_j \in \mathcal{A}$ of alternating percepts and actions is called the t-step history of interaction. At the t^{th} time-step, and given the elapsed history h_{t-1} , the behavior of the agent at step t is given by the map $M_A^{h_{t-1}}(s \in \mathcal{S}) \in distr(\mathcal{A})$,

where $distr(\mathcal{X})$ denotes the set of probability distributions over the set \mathcal{X} . The realized agent's action, given history h_{t-1} , is sampled from the distribution $M_A^{h_{t-1}}(s \in \mathcal{S})$. The environment is specified analogously.

The notion of a reward $\lambda \in \Lambda$ in RL (specifying whether a performed action, or a sequences thereof were 'correct') can be w.l.g. subsumed into the percept space. All the standard figures of merit regarding the learning performance of an agent

are functions of the realized history. They are often convexlinear, enabling statements about average performance. Thus, the interaction history is the fundamental object in RL, and must be maintained in the quantum setting.

In an extension of the above framework to a quantum setting, the percepts/actions are promoted to orthogonal states (percept/action states) of the percept/action Hilbert spaces $\mathcal{H}_{\mathcal{S}} = \operatorname{span}\{|s_i\rangle\}_i$ and $\mathcal{H}_{\mathcal{A}} = \operatorname{span}\{|a_i\rangle\}_i$. The agent, and the environment, contain internal memory: finite (but arbitrarily sized) internal registers R_A and R_E which can store histories, with Hilbert spaces of the form $\mathcal{H}_{\mathcal{A}} \otimes \mathcal{H}_{\mathcal{S}} \otimes \mathcal{H}_{\mathcal{A}} \cdots$. We jointly call the percept/action states, their probabilistic mixtures, and sequences thereof, *classical states*. The interaction is modelled via a common communication register R_C , with associated Hilbert space $\mathcal{H}_C = \{|x\rangle | x \in \mathcal{S} \cup \mathcal{A}\}$ sufficient to represent both actions and percepts whose spaces are mutually orthogonal². The agent (environment) is fully characterized by sequences of completely positive trace preserving (CPTP) maps $\{\mathcal{M}_i^A\}_i$ ($\{\mathcal{M}_i^E\}_i$) acting on the concatenated registers $R_A R_C$ ($R_C R_E$). We assume the three registers are pre-set in a fiducial classical product state. An agent-environment interaction is then defined by the maps which are sequentially applied.

Classical agents (environments) are those whose maps do not generate non-classical states (superpositions of classical states), given classical states of the internal and the communication registers. More generally, an agent and environment have a classical interaction if the joint state of registers $R_A R_C R_E$ is separable w.r.t. the three partitions, and the state of R_C , post-selected on any outcome of any separable measurement of $R_A R_E$ is a classical state, at every stage of the interaction³. To maintain a robust notion of history, we introduce testers, systems which monitor the interaction, and record it in its own memory R_T . Testers we consider are sequences of controlled maps of the form

$$U_{k}^{T}\left(|x\rangle_{R_{C}}\otimes|\psi\rangle_{R_{T}}\right)=|x\rangle_{R_{C}}\otimes U_{k}^{x}\left|\psi\rangle_{R_{T}}$$

where $x \in \mathcal{S} \cup \mathcal{A}$, and $\{U_k^x\}_x$ are unitary maps acting on the register R_T , for all steps k. A tested interaction is shown in Fig. 1. If all the maps of the tester copy⁴ the classical states we call it a *classical tester*. A *sporadic classical tester* allows periods of untested interaction (i.e. some maps are identities).

B. Quantum improvements in learning

Improvements have been demonstrated for a few types of task environments, including epoch-type environments, where the state of the environment is periodically re-set. For example,

²A more general definition of an interaction, where in the spirit of robotics and embodied cognitive sciences we separate the interfaces of the agent and the environment, is provided in [13].

 $^{^3 {\}rm Classical}$ interaction still allows that the internal information processing of the agent and environment includes non-classical states, e.g. an internal quantum computer. However, neither the agent or environment are allowed to output non-classical states, or to be entangled to the communication register R_C .

⁴By 'copy' we mean the map $|x\rangle |\epsilon\rangle \to |x\rangle |x\rangle \ \forall x\in\mathcal{S}\cup\mathcal{A}$, for some fixed state $|\epsilon\rangle$.

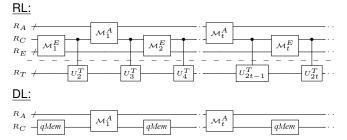


Fig. 1. RL: Tested agent-environment interaction suitable for reinforcement learning. In general, each map of the tester U_k^T acts on a fresh subsystem of the register R_T , which is not under the control of the agent, nor of the environment. The crossed wires represent multiple systems. DL: The simpler setting of standard quantum machine learning, where the environmental map is without internal memory, presented in the same framework. See main text for further detail.

once a chess game finishes, the board is re-set for the next game; similarly, in navigation tasks, e.g. maze navigation, once the agent walker) has found the exit (goal), it is brought back to the initial position, and the task is iterated. Our overall approach is inspired by contrasting environments to oracles which appear in many quantum algorithms. Standard environments do not match the specification of oracles. For instance, the actions of the agent are lost to the environment, and not returned. Nonetheless, for epoch-type environments E, we can define $oracular\ instatiations\ E_q$ of the classically specified environment E, which are unitary. With access to the instatiations, the agent can execute amplitude amplification [15] to obtain an example of a rewarding sequence of actions more efficiently.

Environments E, which can be accessed both in their classical (E) or oracularized instantiations (E_q) , we call controllable environments.

The power to identify winning action sequences alone says little about the learning capacities of the agent. However, an exploration stage, i.e. searching, must precede an exploitation stage, and the correct interplay of these two phases is a well-studied problem in RL [5], [7]⁵. This strongly suggest that (significantly) faster search should be useful for learning. To formalize this intuition we define *luck-favoring settings*. Briefly, a learning model A and environment E are luckfavoring, if a lucky agent (one which by chance alone finds many correct sequences of action during an early stage) outperforms an unlucky agent (one which does not) in E, after this early stage. The performance is measured relative to a chosen figure of merit R. Note that most benchmarking task environments are luck favoring, with standard learning models, relative to the usual figures of merit (e.g. finite-horizon average reward). When we combine luck-favoring settings, with the capacities of quantum agents to explore faster, and the notion of a sporadic classical tester, we obtain the following result (Theorem 1 in [12] simplified):

Main Theorem (informal) Given a classical learning agent A, and a controllable, classically specified epochal task environment E such that (A, E) are luck-favoring relative to some figure of merit R, there exists a quantum agent A^q which outperforms A in R, relative to a chosen sporadic classical tester.

The basic idea of the proof is to have the quantum agent A^q to use quantum untested access to E_q to obtain instances of winning sequences in quadratically reduced time. Given these, A^q will, internally, and by using no interaction steps, 'train' a simulation of A. Eventually, the simulation will get lucky (call this successfully trained simulation A_{lucky}). A^q then relinquishes control to A_{lucky} , and forwards percepts and actions between A_{lucky} and the environment, under a classical tester. Since the setting is luck-favoring, the main claim follows. Although amplitude amplification yields quadratic speed ups in search, in some learning settings this can lead to even an exponential improvement (relative to the task environment size) in learning agent performance. Such exponential separations occur when lucky instances occur exponentially infrequently (e.g. mazes with low connectivity). In this scenario, a quantum agent can find a successful instance with near unit probability while the classical agent still has an exponentially small chance of the same outcome, for a chosen time interval. Thus this can constitute a relevant exponential improvement, for time-limited games.

In previous work, the enhancements were proven by "oracularizing" task environments, to achieve an (abstract) mapping summarized with the following expression

$$|a_1, \dots, a_M\rangle \xrightarrow{E_{oracle}^q} (-1)^{\Lambda(a_1, \dots, a_M)} |a_1, \dots, a_M\rangle,$$
 (1)

where $\Lambda(a_1,\ldots,a_M)$ is 1 iff the sequence of actions a_1,\ldots,a_M results in a reward. This map is then further refined in [12] to achieve similar improvements for a broader class of settings, but the basic idea is the same. Next we show that a similar philosophy, but relying on other oraculizations, can be used for improvements of different flavours, and finally briefly address the conceptual cost of "oraculization" from a practical perspective.

II. QUANTUM IMPROVEMENTS BEYOND STATE EXPLORATION

The basic idea behind quantum improvements described above can be summarized in a sentence. We use quantum access to an (not fully characterized) environment to learn properties which can help us optimize an otherwise classical agent. This leads to the following three-step schema for identifying improvements: a) identify an "indirect property" of interest; b) identify a suitable oraculization (and the corresponding class of environments which are compatible with the oraculization) using which the indirect property will be ascertained and c) prove improvements – i.e. prove that a quantum agent can win overall, although it sacrifices valuable learning steps to ascertain the indirect properties [12]. To exemplify this on the result of the previous section, the relevant property is "path

⁵A related interplay occurs in optimization problems, where a local minimum is often rapidly found, and can be used. However, this opens the possibility that we are missing out on a global minimum.

leading to a reward", the oracle is the blocked oracle marking such paths, and the proof follows from the restriction on luckfavoring environments, and first-success hitting times.

A different way of obtaining improvements considers not just the environment, but also the agent at hand. Most learning models, both in the DL and RL case, come equipped with userdefined settings, called hyperparameters, or metaparameters, in the DL, and RL case, respectively. For example, if one uses neural networks for some task, one first needs to fix an architecture. In the case of RL, if one uses Q-Learning (or SARSA) [7], they must choose the exploration/exploitation transition functions, that is, how quickly the model gets "greedy". The projective simulation (PS) model [10] can have a few parameters including γ (controlling the forgetting speed), η (controlling reward propagation) and so on. The optimal learning settings of these parameters typically depends on the task environment at hand. This is, in fact, unavoidable because of No Free Lunch theorems [16]) which guarantee a universally optimal learner is not possible. In practice, these parameters are set manually by the user. In general this should, and can, be automatized. In recent work, it was shown how these parameters can be learned in the context of the PS, but similar results have been obtained for other models as well (see [17], and references therein). One way to think about this setting is to think of each configuration of parameter settings as fixing some particular learning agent/model A_k . On the other hand, the set of all considered environments is partitioned according to those for which a particular agentconfiguration is optimal (out of the available set $\{A_k\}_k$). This is illustrated in Fig. 2.

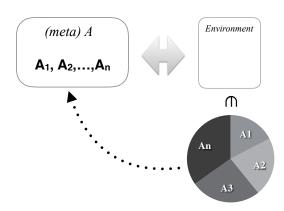


Fig. 2. Metalearning: The learning agent, over time, identifies the optimal configuration of its own metaparameters in order to optimize performance.

Using the same oraculization approach as previously, we can now show how metalearning itself can be quantum-enhanced. The exact derivation of a possible approach is involved and long but the basic ideas are given next. First, note that, by taking a birds-eye perspective view of the scenario, any learning model with hyper/meta-parameters can be represented as a controlled unitary, controlled by a register specifying the values of the metaparameters (represented in the state $|\mathbf{k}\rangle_{m.m.p.}$, where m.m.p stands for "model metaparameters"). With \mathbf{k} we denote any collection of such parameters. The

"model unitary" is acting on the current memory of the agent (containing so-called model parameters, abbreviated m.p.), and the communication channel. As usual, we represent the environment in its purified picture, which simply means no information is lost. A learning agent without metalearning would then, through interaction with the environment, realize a joint agent-channel-environment state of the form

$$|non - M.A.\rangle = |\mathbf{k}\rangle_{m,m,p,} |mem_A\rangle_{m,p,} |a\rangle_C |state\rangle_E,$$

where mem_A is the memory of the agent, a is the current action (assuming it is the agent's move), and state is the purified state of the environment. The expression above assumes both the agent and environment are described by deterministic programs. In the case either is randomized, the systems of the agent and environment would be entangled, with an extra system, specifying the random choices, but we ignore this for the moment. Note that since everything is unitary, this overall state must determine the entire history of interaction. However k does not depend on it. In order to metalearn the parameters, the standard approach is to have the agent monitor its own performance. Such a pre-metalearning agent would realize the joint state of the form

$$\left|p.M.A.\right\rangle = \left|\mathbf{k}\right\rangle_{m.m.p.} \left|eval(\mathbf{k})\right\rangle \left|mem_A\right\rangle_{m.p.} \left|a\right\rangle_C \left|state\right\rangle_E,$$

where $eval(\mathbf{k})$ is some objective evaluation of the performance (the value depends on the environment, agent and \mathbf{k}), under settings \mathbf{k} , in E.

Now, in classical metalearning, a full metalearning agent must somehow optimize $eval(\mathbf{k})$ in \mathbf{k} , and this is is usually done by gradient descent. The idea here, naturally, is to use quantum processing to optimize this parameter. There are many proposals for quantum optimization (including recent [19]), applicable in various cases, but many assume a special form of the function to be optimized. Since here we wish to maintain the presentation model-independent (albeit, at the cost of its efficiency) we approach this, for now, using the generic black-box, discrete optimization techniques such as [18]. Note, in doing so we have sacrificed the potential for exponential improvements. Such approaches typically "Groverize" the process, and we proceed likewise here: The agent would first initialize its metaparameter register in a uniform superposition, so proportional to $\sum_{\mathbf{k}}$, leading to the superposed state

$$|M.A_q.\rangle \propto \sum_{\mathbf{k}} |\mathbf{k}\rangle_{m.m.p.} |eval(\mathbf{k})\rangle |mem_A\rangle_{m.p.} |a\rangle_C |state\rangle_E.$$

Note, so far we did not need to adjust the environment in any way, and if the agent were to now observe its evaluation register, it would collapse its history to match one particular choice of \mathbf{k} (corresponding to $eval(\mathbf{k})^6$). If it is lucky, or if this is repeated sufficiently many times, it would find the optimal choice, and this is eqivalent to a classical "random guess" agent. If the environment is epochal (as in [12]),

⁶The metaparameter register may still be in superposition if multiple choices lead to the same value.

and controllable to the extent that we can reverse it, we can now amplitude amplify the history-branch corresponding to the best k.7 For this process to be feasible, it is not necessary to know the exact value of the optimal performance $eval(\mathbf{k})$, but just a reasonable range. In this case, one can additionally use a binary search to pinpoint the optimal value while yielding only a logarithmic overhead, relative to the setting where the optimal value is known [18]). The above method is optimal, and yields a strict quadratic improvement, when the mapping from metaparameters to performance has absolutely no structure. Often, we are facing models at the opposite extreme where the parameter space has plenty of structure. For example, in [17], the meta-parameter γ – the "forgetting" parameter, but the details are not relevant here - of the PS model is, is optimized. The analysis of this work also suggests that, in many environments, the abstract mapping, which maps the parameter value to and averaged performance $[0,1] \stackrel{eval}{\rightarrow} \mathbb{R}^+$, is an unimodal function – it is increasing to some value γ_{opt} , and decreasing afterwards. This case is easy for classical learning, as one can perform a binary search. However, even here there are better options, and, in terms of expected time, probably the best option is using the quantum Bayesian adaptive learner [23], leading to a three-fold improvement (note, here improvements beyond multiplicative are not possible). The intermediate regimes are more interesting, specifically when we deal with many continuous parameters. In this case one can, for instance, use "quantum-improved" gradient descent methods [20], or use "completely quantum" gradient descent techniques [19]. Once the learning model is fixed, it is likely that even more efficient methods become available, but this is beyond the scope of this paper.

III. COST OF ORACULIZATION

Most of the analyses we had done simply assume access to controllable environments, and consequently, to oracularized instantiations of the task environment. This is a reasonable assumption if the environments are constructed. If not, sometimes oracularized instantiations can nonetheless be effectively realized. In particular, this is possible if we relax the model of interaction, and grant additional powers to the agent: register hijacking, and register scavenging [13]. Hijacking pertains to the option that the agent has access to the register (memory) of the environment (but not to the specification of the environment it is actually learning), or to the purifying systems of the environment (which keep track of the history). Here we show that these two options suffice implement the useful oraculizations in a (near) black-box fashion (i.e. without needing the specificities of the implementation).

For illustrative purposes, consider deterministic environments with one reward. For this particular case, we need to construct the map in Eq. 1.

We will give an explicit construction below, simplifying the steps given in [13], in which we assume the environment realizes a certain information transfer using specific gates – in reality, whatever the environment does corresponds to *some* choice of maps, which are isomorphic to the concrete choices we make below. In this sense, our choices are without the loss of generality. The steps of the construction are thus as follows.

First using minimal register hijacking, the agent can, effectively, realize the *phase kick-back map*

$$|a_1, \dots, a_M\rangle_A |\epsilon, \dots, \epsilon\rangle_E |\epsilon, \phi^-\rangle_E$$

$$\mapsto (-1)^R |a_1, \dots, a_M\rangle_E |s_2, \dots, s_M\rangle_A |s_{M+1}, \phi^-\rangle_A,$$

where the states $|\epsilon, \phi^-\rangle$ and $|s_{M+1}, \phi^-\rangle$ will be specified later. The map above is just a representation of M steps of

The map above is just a representation of M steps of a standard interaction, where the agent committed to the action sequence a_1,\ldots,a_M , to which the environment responded with the percepts $|s_2,\ldots,s_M\rangle\,|s_{M+1},\phi^-\rangle$ (hence, the agent no longer has access to the actions the environment stored/remembered.). The register initialized in the state $|\epsilon,\phi^-\rangle$ is actually in the possession of the environment, and it should contain the memory slot where the environment will set the last percept, and the reward value. We can assume the reward value is flipped using a two-level unitary: Pauli-X (σ_x) , flipping between $|0\rangle$ to $|1\rangle$. Here, the agent intervenes via hijacking, and imprints the state $|\phi^-\rangle$, which is the -1 eigenstate of σ_x . This will cause a phase kick-back.

Note that the map above is not yet the needed oracle as the action register is still entangled to the percept register. Using a similar construction, the agent can also realize the raw map for the percepts (with no reward status):

$$|a_{1}, \dots, a_{M}\rangle_{A} |\epsilon, \dots, \epsilon\rangle_{E} |\epsilon, \phi^{+}\rangle_{E}$$

$$\mapsto |a_{1}, \dots, a_{M}\rangle_{E} |s_{2}, \dots, s_{M}\rangle_{A} |s_{M+1}, \phi^{+}\rangle_{A}.$$

The final ingredient of the construction is to note that any environment can be described using maps of which are selfreversible.

Since the environment is classically specified, all the maps it realizes can be represented in a particular controlled form – all classical computations also have this property. More precisely, the environments (can be assumed to) apply the map E_t specified with

$$E_t | a_1, \dots, a_t \rangle \otimes | \epsilon \rangle = | a_1, \dots, a_t \rangle \otimes E^{a_1, \dots, a_t} | \epsilon \rangle,$$
 (2)

at each time-step t. These maps E^{a_1,\dots,a_t} only rotate the fiducial (empty) state $|\epsilon\rangle$ to $|s_{t+1}(a_1,\dots,a_t)\rangle$ (indicating also that s_{t+1} depends on the previous actions). As clarified, the reward-setting unitaries act non-trivial only on a two-level subspace, but so are the maps E^{a_1,\dots,a_t} , despite operating on an $|\mathcal{S}|$ - dimensional Hilbert space. But then, they can be chosen such that each E^{a_1,\dots,a_t} is Hermitian (and non-trivial only on a two-dimensional subspace), thus self-inverse.

Note, the capacity that the environment can be reversed implies either a deterministic environment, or that, in the register hijacking phase [13], the agent can temporarily access the register purifying the random choices of the environment. However, it is interesting to consider purely deterministic environments as well, where the agent need only have access to its own randomness-inducing purifying register. This we can assume to be possible without loss of generality.

The total process of oraculization is thus given by the following steps: a) The agent realizes the phase kick-back map, using hijacking; b) again, using hijacking, it implants the percept-containing systems (which it collected previously using scavenging) into the environment, at the correct memory location; c) the agent realizes the raw map for percepts. Given that the maps are all self-reversible (involutions), in total this realizes the oracular instantiation of E_{oracle} . The total cost of interactions (number of interaction steps) with the environment we used was 2M, as two full games were played.

While we naturally do not claim that the above process can be done with real macroscopic environments, it does show that there is no problem in principle, and that the definitional parts of the environment (i.e. what behavior is rewarded and so on), need not be touched to achieve oraculization. If the environment is say a computational environment of a future quantum network (which maintains the purity of all systems), the above is fully possible. It also becomes a possiblity in the case of nano-scale robots in future quantum experiments, where the environment is manifestly quantum and exquisitely controled, as mentioned in [12]. Critically, these processes are fully realizable to speed-up model-based learning, where the learner builds an internal representation of the environment, which can be fully quantum and controllable⁸.

IV. CONCLUSION

The recent successes in classical AI and quantum machine learning point towards the possibilities of new, perhaps revolutionary, quantum-enhanced smart technologies. However, most such smart technologies, aside from performing extensive data analysis, must still learn from interaction. This opens a "quantization bottleneck" which, for instance, prevents us from seriously discussing actual quantum artificial intelligence, (thus beyond clustering and classification). Progress in quantum reinforcement learning may help bridge that gap. In this work, we have addressed the recent efforts in this field, and expanded the results for quantum improvements in learning quality, by showing how the process of learning how to learn – metalearning – can itself be quantum-enhanced. This yields quantum-enhanced learning agents which, generically, outperform their classical counterparts. As a side result, we have presented a succinct construction which shows how oracular instantiations of task environments, necessary for improvements, can be realized in a black-box fashion. This latter result may be relevant in future implementations of quantum reinforcement learning.

ACKNOWLEDGMENT

VD and HJB acknowledge the support by the Austrian Science Fund (FWF) through the SFB FoQuS F 4012, and the Templeton World Charity Foundation grant TWCF0078/AB46.

REFERENCES

- [1] Wittek, P. Quantum Machine Learning: What Quantum Computing Means to Data Mining (Academic Press, 2014).
- [2] Lloyd, S., Mohseni, M. and Rebentrost, P. Quantum algorithms for supervised and unsupervised machine learning. ArXiv:1307.0411 (2013).
- [3] Aïmeur, E., Brassard, G. and Gambs, S. Quantum speed-up for unsupervised learning. *Machine Learning* 90, 261–287 (2013).
- [4] Rebentrost, P., Mohseni, M, and Lloyd, S. Quantum Support Vector Machine for Big Data Classification. *Phys. Rev. Lett.* 113, 130503 (2014).
- [5] Sutton, R. S. and Barto, A. G. Reinforcement learning: An introduction (MIT Press, Cambridge Massachusetts, 1998), first edn.
- [6] Silver, D, et. al. Mastering the game of Go with deep neural networks and tree search Nature 529, 484–503 (2016).
- [7] Russel, S. J. and Norvig, P. Artificial intelligence A modern approach (Prentice Hall, New Jersey, 2003), second edition edn.
- [8] Dong, D., Chunlin, C. and Zonghai, C. Quantum Reinforcement Learning. Advances in Natural Computation Lecture Notes in Computer Science 3611 686-689 (2005).
- [9] Paparo, G. D., Dunjko, V., Makmal, A., Martin-Delgado, M. A. and Briegel, H. J. Quantum speedup for active learning agents. *Phys. Rev.* X 4, 031002 (2014).
- [10] Briegel, H. J. & De las Cuevas, G. Projective simulation for artificial intelligence. Sci. Rep. 2 (2012).
- [11] Giovannetti, V., Lloyd, S. and Maccone, L. Quantum Random Access Memory. Phys. Rev. Lett. 100, 160501 (2008).
- [12] Dunjko, V., Taylor, J. M. & Briegel, H. J. Quantum-Enhanced Machine Learning. Phys. Rev. Lett. 117, 130501 (2016).
- [13] Dunjko, V., Taylor, J. M. & Briegel, H. J. Framework for learning agents in quantum environments. arXiv:1507.08482 (2015).
- [14] Chiribella, G., D'Ariano, G. M. & Perinotti, P. Quantum Circuit Architecture. Phys. Rev. Lett. 101, 060401 (2008).
- [15] Brassard, G., Hoyer, P., Mosca, M. and Tapp, A. Quantum Amplitude Amplification and Estimation arXiv:quant-ph/0005055 (2000).
- [16] Wolpert, D. H. The lack of a priori distinctions between learning algorithms. *Neural Comput.* 8, 1341–1390 (1996).
- [17] Makmal, A, Melnikov, A. A., and Briegel, H. J. Meta-learning within Projective Simulation. *IEEE Access* 4, 2110 (2016).
- [18] Durr, C., Hoyer, P. A Quantum Algorithm for Finding the Minimum. arXiv:quant-ph/9607014 (1996).
- [19] Rebentrost, P., Shuld, M., Petruccione, F. and Lloyd, S. Quantum gradient descent and Newton's method for constrained polynomial optimization. arXiv:1612.01789v1 (2016).
- [20] Jordan, S. P. Fast Quantum Algorithm for Numerical Gradient Estimation. Phys. Rev. Lett. 95, 050501 (2005).
- [21] Schuld, M., Sinayskiy, I., Petruccione, F. Prediction by linear regression on a quantum computer. *Phys. Rev. A* 94, 022342 (2016).
- [22] Lloyd, S., Garnerone, S. & Zanardi, P. Quantum algorithms for topological and geometric analysis of data. *Nat. Comm.* 7, 10138 (2016).
- [23] Ben-Or, M. and Hassidim, A.,he Bayesian Learner is Optimal for Noisy Binary Search (and Pretty Good for Quantum As Well). Proceedings of the 2008 49th Annual IEEE Symposium on Foundations of Computer Science 221-230 (2008).

⁸Indeed, in this case it may be possible to build the oracular variant directly, however the described procedure shows how this can be done in a black-box fashion, i.e. without caring about the details of implementation.