

# Approximate Dynamic Programming using Fluid and Diffusion Approximations with Applications to Power Management

Wei Chen\*, Dayu Huang\*, Ankur Kulkarni†, Jayakrishnan Unnikrishnan\*, Quanyan Zhu\*,  
Prashant Mehta†, Sean Meyn\*, and Adam Wierman§

\*Dept. of ECE and CSL, UIUC, Urbana, IL 61801, U.S.A.

†Dept. of MSE and CSL, UIUC, 1206 W. Green Street, Urbana, IL 61801

‡Dept. of Industrial and Enterprise Sys Engg (IESE), UIUC, 104 S. Mathews Avenue, Urbana, IL 61801

§Dept. of CS, California Inst. of Tech.. 1200 E. California Boulevard. Pasadena, CA, 91125.

**Abstract**—TD learning and its refinements are powerful tools for approximating the solution to dynamic programming problems. However, the techniques provide the approximate solution only within a prescribed finite-dimensional function class. Thus, the question that always arises is *how should the function class be chosen?* The goal of this paper is to propose an approach for TD learning based on choosing the function class using the solutions to associated fluid and diffusion approximations. In order to illustrate this new approach, the paper focuses on an application to *dynamic speed scaling* for power management.

**Keywords:** Nonlinear control, adaptive control, machine learning, optimal stochastic control, dynamic speed scaling.

**AMS subject classifications:** Primary: 93E35, 49J15, 93C40  
Secondary: 65C05, 93E20 68M20

## I. INTRODUCTION

Stochastic dynamic programming and, specifically, controlled Markov chain models (MDPs) have become central tools for evaluating and designing communication, computer, and network applications. These tools have grown in popularity as computing power has increased; however, even with increasing computing power, it is often impossible to attain exact solutions. This is due to the so-called “curse of dimensionality”, which refers to the fact that the complexity of dynamic programming equations often grows exponentially with the dimension of the underlying state space.

However, the “curse of dimensionality” is slowly dissolving in the face of approximation techniques such as Q-learning and TD-learning [26], [6], [10]. These techniques are designed to approximate a solution to a dynamic programming equation within a prescribed finite-dimensional function class. A key determinant of the success of these techniques is the selection of this function class. The question of how to select an appropriate basis has been considered in specific contexts, e.g. [27], [16]. However, despite the progress so far, determining the appropriate function class for these techniques is still more of an art than a science.

The goal of this paper is to illustrate that a useful function class can be attained by solving the dynamic programming equation for a highly idealized approximate model. Specifically, a useful function class is obtained by first constructing a fluid or diffusion approximation of the MDP model, and solving the corresponding dynamic programming equation for the simpler system.

In the special case of network scheduling and routing, it is known that the dynamic programming equations for the continuous-time model are closely related to the corresponding equations for the discrete-time model. This relationship has been developed by one of the authors in [19], [12], [20] and the monograph [21]. Moreover, the fluid value function has been used as part of a basis in the approximate dynamic programming approaches of [28], [22]. In this paper we demonstrate that the solution to the dynamic programming equations for the fluid, diffusion, and discrete-time models are closely related in more general classes of models.

In order to provide a concrete illustration of the proposed approximation techniques, the paper considers an example of a stochastic control problem from the area of power management in computer systems. Specifically, an important tradeoff in modern computer system design is between reducing energy usage and maintaining good performance (small delays). To this end, an important technique is *dynamic speed scaling* [4], [36], [34], [13], which dynamically adjusts the processing speed in response to changes in the workload — reducing (increasing) the speed in times when the workload is small (large). Dynamic speed scaling is now common in many chip designs, e.g. [2], [1], and network environments, e.g. switch fabrics [17], wireless communication [32], [8], and TCP offload engines [23]. Further, dynamic speed scaling has been the focus of a growing body of analytic research [3], [5], [11], [31]. Sec. III provides the details about the speed scaling model and reviews related literature.

For purposes of this paper, dynamic speed scaling is simply a stochastic control problem – a single server queue with a controllable service rate – and the goal is to understand how to control the service rate in order to minimize the total cost, which is a weighted sum of the energy cost and the delay cost. In this context, this paper will illustrate how to use the solutions of the fluid and diffusion models in order to apply TD learning to determine an approximately optimal policy for control. Fluid and diffusion models for the dynamic speed scaling problem are analyzed in Sec. IV. The results of applying TD learning to the speed scaling problem are illustrated in Sec. V. These results highlight the usefulness of the fluid and diffusion solutions for TD learning.

Fig. 1 illustrates the results obtained using TD learning in

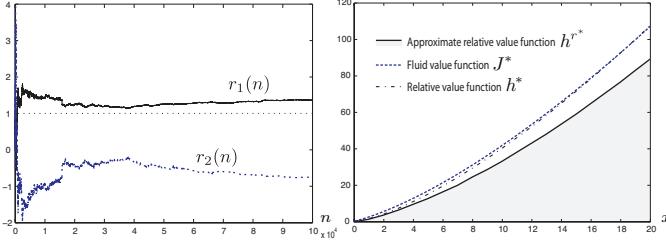


Fig. 1: Simulation results for the dynamic speed scale model with quadratic cost. The plot on the left shows estimates of the coefficients in the optimal approximation of  $h^*$  using the basis obtained from the fluid and diffusion models (see (36)). In the plot on the right the final approximation  $h^{r^*}$  is compared to the fluid value function and the relative value function.

this application. The three plots compared in the figure are the fluid value function  $J^*$  appearing in the Total Cost Optimality Equation (5), the relative value function  $h^*$  appearing in the Average Cost Optimality Equation (3), and the approximate value function obtained from the TD learning algorithm. The basis obtained from analysis of the fluid and diffusion models results in a remarkably tight approximation of  $h^*$ .

Although the paper focuses in large part on the application of TD learning to the dynamic speed scaling problem, the approach presented in this paper is general. The use of fluid and diffusion approximations to provide an appropriate basis for TD learning is broadly applicable to a wide variety of stochastic control problems.

## II. PRELIMINARIES

### A. Markov Decision Processes (MDPs)

In this paper we will consider the following general MDP model. Let  $\mathcal{X} = \mathbb{R}_+^\ell$  denote the state space for the model. The action space is denoted  $\mathcal{U}$ . In addition there is an i.i.d. process  $\mathbf{W}$  evolving on  $\mathbb{R}^w$  that represents a disturbance process. For a given initial condition  $X(0) \in \mathcal{X}$ , and a sequence  $\mathbf{U}$  evolving on  $\mathcal{U}$ , the state process  $\mathbf{X}$  evolves according to the recursion,

$$X(t+1) = X(t) + f(X(t), U(t), W(t+1)), \quad t \geq 0. \quad (1)$$

We restrict to inputs that are defined by a (possibly randomized) stationary policy. This defines a Markov Decision Process (MDP) with controlled transition law

$$P_u(x, A) := \mathbb{P}\{x + f(x, u, W(1)) \in A\}, \quad A \in \mathcal{B}(\mathcal{X}).$$

We let  $\mathcal{D}_u$  denote the generator in discrete time. For any function  $h: \mathbb{R} \rightarrow \mathbb{R}$ ,

$$\mathcal{D}_u h(x) := \mathbb{E}[h(X(t+1)) - h(X(t)) | X(t) = x, U(t) = u] \quad (2)$$

A cost function  $c: \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}_+$  is given, and our goal is to find an optimal control based on this cost function. We focus on the average cost problem, with associated Average Cost Optimality Equation (ACOE):

$$\min_u (c(x, u) + \mathcal{D}_u h^*(x)) = \eta^* \quad (3)$$

The ACOE is a fixed point equation in the *relative value function*  $h^*$ , and the optimal cost for the MDP  $\eta^*$ .

### B. The fluid and diffusion models

The fluid model associated with the MDP model is defined by the following mean flow equations,

$$\frac{d}{dt} x(t) = \bar{f}(x(t), u(t)), \quad x(0) \in \mathcal{X},$$

where  $u$  evolves on  $\mathcal{U}$ , and  $\bar{f}(x, u) := \mathbb{E}[f(x, u, W(1))]$ . The generator for the fluid model is defined similarly. Given  $u(0) = u$ ,  $x(0) = x$ ,

$$\mathcal{D}_u^F h(x) = \left. \frac{d}{dt} h(x(t)) \right|_{t=0} = \nabla h(x) \cdot \bar{f}(x, u). \quad (4)$$

The associated Total Cost Optimality Equation (TCOE) is

$$\min_u (c(x, u) + \mathcal{D}_u^F J^*(x)) = 0 \quad (5)$$

It is solved with the value function,

$$J^*(x) = \inf_u \int_0^\infty c(x(t), u(t)) dt, \quad x(0) = x \in \mathcal{X}, \quad (6)$$

provided  $J^*$  is finite valued, which requires assumptions on the cost and dynamics. Under these assumptions the optimal policy is any minimizer,

$$\phi^{F*}(x) \in \arg \min_u (c(x, u) + \mathcal{D}_u^F J^*(x)) \quad (7)$$

In many models, e.g. queueing networks, the applicability of the fluid model can be justified through a scaling argument similar to the following: For a large initial condition, and over a long time horizon, the sample paths of the stochastic model can be approximated by a solution to the fluid model equations. Based on this approach, techniques have been developed that provide easily verified stability conditions for stochastic networks based on the analysis of the fluid model.

Similar scaling arguments can be used to show that  $h^*$  is approximated by  $J^*$ . For history and further results see [21]. However, this approach *fails* in the example considered in the current paper because the input is not bounded. Thus we need a different motivation for considering the fluid model. Here, motivation for approximate models comes from a Taylor series expansion. In particular, if the fluid value function  $J^*$  is smooth then we have the approximation,

$$\begin{aligned} \mathcal{D}_u J^*(x) &\approx \mathbb{E}_{x,u} [\nabla J^*(X(0))(X(1) - X(0))] \\ &= \nabla J^*(x) \bar{f}(x, u) \end{aligned} \quad (8)$$

where the subscript indicates expectation conditional on  $X(0) = x$ ,  $U(0) = u$ . That is,  $\mathcal{D}_u J^* \approx \mathcal{D}_u^F J^*$ , where the approximation depends on the smoothness of the function  $J^*$ . In the example treated in Sec. IV-A we obtain precise error bounds which illustrate that  $J^*$  almost solves the ACOE for the stochastic model.

A diffusion model is obtained similarly. We again choose its dynamics to reflect the behavior of the discrete-time model. To capture the state space constraint we opt for a reflected diffusion, defined by the Ito equation:

$$dX(t) = \bar{f}(X(t), U(t))dt + \sigma(U(t))dN(t) + dI(t), \quad (9)$$

where the process  $N$  is a standard Brownian motion on  $\mathbb{R}^\ell$  and  $I$  is a reflection process. That is, for each  $1 \leq i \leq \ell$ ,

the process  $I_i$  is non-decreasing and is minimal subject to the constraint that  $X_i(t) \geq 0$  for each  $t$  and each  $i$ . This is captured through the sample path constraint,

$$\int_0^\infty X_i(t) dI_i(t) = 0, \quad 1 \leq i \leq \ell \quad (10)$$

For more on reflected diffusions see [30].

In the dynamic speed scaling model we find that the ACOE associated with the diffusion model is approximately solved by a perturbation of  $J^*$ . This depends, of course, on the choice of the variance term  $\sigma^2(u)$  in (9). In Sec. IV-B we argue that this should be chosen based on a second-order Taylor-series approximation of the ACOE for the primary discrete-time model, much like the first-order Taylor series approximation (8) that motivated the generator for the fluid model.

### C. TD learning

TD learning is a technique for approximating value functions of MDPs within a linearly parameterized class.

Specifically, we define  $\{\psi_i : 1 \leq i \leq d\}$  as real-valued functions on  $X$  and we let  $h^r = \sum r_i \psi_i$  or, with  $\psi: X \rightarrow \mathbb{R}^d$  the vector of basis functions,  $h^r = r^T \psi$ . Suppose that a stationary policy  $\phi$  is applied to the MDP model, and that the resulting Markov chain is ergodic with stationary marginal  $\pi$ . Let  $h$  denote the solution to Poisson's equation  $P_\phi h = h - c_\phi + \eta_\phi$  where  $P_\phi(x, dy) = P_{\phi(x)}(x, dy)$  is the resulting transition law for the chain,  $c_\phi(x) = c(x, \phi(x))$  is the cost as a function of state for this policy, and  $\eta_\phi$  is the average cost. TD learning then takes the mean-square error criterion:

$$\frac{1}{2} \mathbb{E}_\pi[(h(X(0)) - h^r(X(0)))^2] := \frac{1}{2} \int (h(x) - h^r(x))^2 \pi(dx).$$

Hence the optimal parameter satisfies the fixed point equation,

$$\mathbb{E}_\pi[(h(X(0)) - h^r(X(0)))\psi(X(0))] = 0. \quad (11)$$

In the rest of this section we assume that the control is fixed to be  $\phi(x)$ . We use  $c(x)$  to denote the cost function  $c_\phi(x)$  and  $\mathbb{E}$  to denote the expectation under this stationary policy.

The TD and LSTD learning algorithms are techniques for computing the optimal parameter. We refer the reader to Chapter 11 of [21] for details of the LSTD learning algorithm used in the numerical results described in this paper and provide only a high-level description of the LSTD algorithm here.

When the parameterization is linear then (11) implies that the optimal parameter can be expressed

$$\begin{aligned} r^* &= \Sigma^{-1} z \quad \text{with } \Sigma = \mathbb{E}_\pi[\psi(X(0))\psi(X(0))^T] \\ z &= \mathbb{E}_\pi[\psi(X(0))h(X(0))]. \end{aligned} \quad (12)$$

The matrix  $\Sigma$  can be estimated using sample path averages of  $\{\psi(X(t))\psi(X(t))^T\}$ . The same is true for  $z$ , following a transformation.

This transformation requires some machinery: First, it is known that the solution to Poisson's equation can be expressed  $h = Zc$ , where  $Z$  is the *fundamental kernel*. Under appropriate

conditions on the Markov chain this is expressed as the conditional expectation,

$$Zc(x) = \mathbb{E}\left[\sum_{t=0}^{\tau_{x^*}-1} [c(X(t)) - \eta] \mid X(0) = x\right]$$

where  $x^*$  is any fixed state with non-zero steady-state probability, and  $\tau_{x^*} \geq 1$  is the first entrance time. The representation of  $z$  is obtained in a Hilbert space setting, based on the adjoint of  $Z$ . Let  $L_2$  denote the usual Hilbert space of square-integrable functions, with inner product  $\langle f, g \rangle = \mathbb{E}[f(X(0))g(X(0))]$ ,  $f, g \in L_2$ . Letting  $Z^\dagger$  denote the adjoint of  $Z$  gives,

$$z_i = \langle Zc, \psi_i \rangle = \langle c, Z^\dagger \psi_i \rangle$$

This representation is useful for estimation because the adjoint can be expressed in terms of the stationary process on the two-sided time axis,  $\{X(t) : -\infty < t < \infty\}$ . Let  $\tau_{x^*}^-$  denote the *last time* prior to  $t = 0$  that  $x^*$  was visited. Then, for any  $f \in L_2$  with mean  $\eta_f$  we have,

$$Z^\dagger f(x) = \mathbb{E}\left[\sum_{\tau_{x^*}^- < t \leq 0} [f(X(t)) - \eta_f] \mid X(0) = x\right]$$

provided the expectation exists and is finite-valued.

To estimate  $z$  we define the sequence of *eligibility vectors*,

$$\varphi(t+1) = \varphi(t) + \mathbb{I}\{X(t) \neq x^*\}(\psi(X(t)) - \eta_\psi(t))$$

where  $\varphi(0) = \psi(X(0))$ , and  $\eta_\psi(t)$  the sample mean of  $\psi$ . We then define,

$$\Sigma_T = \frac{1}{T} \sum_{t=1}^T \psi(X(t))\psi^T(X(t)), \quad z_T = \frac{1}{T} \sum_{t=1}^T c(X(t))\varphi(t)$$

The LSTD learning algorithm for average cost defines estimates of  $r^*$  in (12) via,

$$r_T = \Sigma_T^{-1} z_T$$

This is consistent provided  $\psi$  and  $h$  are square integrable.

### III. POWER MANAGEMENT VIA SPEED SCALING

This paper proposes a general technique for choosing a basis for TD learning. However, in order to ground the proposed approach, we focus on a specific example of a stochastic control problem that is of particular importance to modern computer system design: *dynamic speed scaling*.

Dynamic speed scaling is an increasingly common approach to power management in computer system design. The goal is to control the processing speed so as to optimally balance energy and delay costs – reducing (increasing) the speed in times when the workload is small (large).

We model the dynamic speed scaling problem as a single server queue with controllable service rate. Specifically, we assume that jobs arrive to a single processor and are processed at a rate determined by the current power. The primary model is described in discrete time: For each  $t = 0, 1, 2, \dots$  we let  $A(t)$  denote the job arrivals in this time slot,  $Q(t)$  the number of jobs awaiting service, and  $U(t)$  the number of services. It

is assumed that  $\mathbf{A}$  is i.i.d. Hence the MDP model is described as the controlled random walk,

$$Q(t+1) = Q(t) - U(t) + A(t+1), \quad t \geq 0. \quad (13)$$

This is an MDP model of the form (1) with  $\mathbf{X} \equiv \mathbf{Q}$ . The cost function we consider balances the cost of delay with the energy cost associated with the processing speed:

$$c(x, u) = x + \beta \mathcal{P}(u), \quad (14)$$

where  $\mathcal{P}$  denotes the power required as a function of the speed  $u$ , and  $\beta > 0$ . This form of cost function is common in the literature, e.g., [11], [5], [31].

The remaining piece of the model is to define the form of  $\mathcal{P}$  — an appropriate form is highly application dependent. In this paper, we consider two particular application areas where speed scaling is an important approach: processor design and wireless transmission.

In the domain of processor design, prior literature has typically assumed  $\mathcal{P}$  is a polynomial, specifically a cubic. That is because the dynamic power of CMOS is proportional to  $V^2 f$ , where  $V$  is the supply voltage and  $f$  is the clock frequency [15]. Operating at a higher frequency requires dynamic voltage scaling (DVS) to a higher voltage, nominally with  $V \propto f$ , yielding a cubic relationship. However, recent work, e.g. [31], has found that the dynamic power usage of real chips is well modeled by a polynomial scaling of speed to power, but this polynomial is closer to quadratic. Thus, in this case we take:

$$\mathcal{P}(u) \propto u^\varrho \quad (15)$$

where  $\varrho > 1$ , but we often focus on the case of  $\varrho = 2$ .

In the case of wireless transmissions, the form of  $\mathcal{P}(u)$  differs significantly. In particular, considering an additive white Gaussian noise model [32] gives, for some  $\kappa > 0$ ,

$$\mathcal{P}(u) \propto e^{\kappa u} \quad (16)$$

There is a large literature on the dynamic speed scaling problem, beginning with Yao et al. [33]. Much of the work focuses on models with either fixed energy budgets [25], [7], [35] or job completion deadlines [24], [3]. In the case where the performance metric is the weighted sum of energy and delay costs (as in the current paper), a majority of the research is in a deterministic, worst-case setting [3], [5]. Most closely related to the current paper are [11], [31], which consider the MDP described above. However, neither of these papers consider either the fluid or diffusion approximations of the speed scaling model; nor do they discuss the application of TD learning.

#### IV. APPROXIMATE MODELS

In this section we study the fluid and diffusion approximations of the speed scaling model described in (13). The solutions to these approximate models will later serve as the basis for applying TD learning to determine an approximately optimal control of the speeds.

##### A. The fluid model

The fluid model corresponding to the speed scaling model (13) is given by:

$$\frac{d}{dt} q(t) = -u(t) + \alpha, \quad (17)$$

where  $\alpha$  is the mean of  $A(t)$ , and the control  $u(t)$  and buffer contents  $q(t)$  are assumed to be non-negative valued.

It is assumed here that the cost function vanishes at the equilibrium  $q(t) = 0$ ,  $u(t) = \alpha$ . In this case the total cost  $J^*$  defined in (6) is finite for each  $x$ . The infimum in (6) is over all feasible  $u$ . Feasibility means that  $u(t) \geq 0$  for each  $t$ , and the resulting state trajectory  $q$  is also non-negative valued. In this section we consider two classes of normalized cost functions,

$$\begin{aligned} \text{Polynomial cost } c(x, u) &= x + \beta([u - \alpha]_+)^{\varrho} \\ \text{Exponential cost } c(x, u) &= x + \beta[e^{\kappa u} - e^{\kappa \alpha}]_+ \end{aligned} \quad (18)$$

where  $[\cdot]_+ = \max(0, \cdot)$ , and the parameters  $\beta, \kappa, \varrho$  are positive. The normalization is used to ensure that  $c(0, \alpha) = 0$ . Observe that the cost is also zero for  $u < \alpha$  when  $x = 0$ . However, it can be shown that the  $u$  that achieves the infimum in (6) is never less than  $\alpha$ .

We now return to (8) to show that the fluid value function provides a useful approximation to the solution to the average cost optimality equations. We construct a cost function  $c^\circ$  that approximates  $c$ , along with a constant  $\eta^\circ > 0$  such that  $J^*$  satisfies the ACOE for this cost function:

$$\min_{0 \leq u \leq x} \{c^\circ(x, u) + P_u J^*(x)\} = J^*(x) + \eta^\circ. \quad (19)$$

This construction is based on the two error functions,

$$\begin{aligned} \mathcal{E}(x, u) &= c(x, u) - J^*(x) + P_u J^*(x) \\ \underline{\mathcal{E}}(x) &= \min_{0 \leq u \leq x} \mathcal{E}(x, u) \end{aligned} \quad (20)$$

The constant  $\eta^\circ \in \mathbb{R}_+$  is arbitrary, and the perturbation of the cost function is defined as

$$c^\circ(x, u) = c(x, u) - \underline{\mathcal{E}}(x) + \eta^\circ$$

Based on the definition of  $\underline{\mathcal{E}}$ , we conclude that (19) is satisfied. To demonstrate the utility of this construction it remains to obtain bounds on the difference between  $c$  and  $c^\circ$ .

We begin with some structural results for the fluid value function. Proofs are omitted due to lack of space. Note that part (ii) is obtained from bounds on the “Lambert  $W$  function” [14].

**Proposition 1.** *For any of the cost functions defined in (18), the fluid value function  $J^*$  is increasing, convex, and its second derivative  $\nabla^2 J^*$  is non-increasing. Moreover,*

(i) *For polynomial cost the value function and optimal policy are given by, respectively,*

$$J^*(x) = x^{\frac{2\varrho-1}{\varrho}} \frac{\varrho}{2\varrho-1} \left( \frac{1}{\beta(\varrho-1)} \right)^{\frac{\varrho-1}{\varrho}} \quad (21)$$

$$\phi^{F*}(x) = \left( \frac{x}{\beta(\varrho-1)} \right)^{1/\varrho} + \alpha, \quad x \in \mathbb{R}_+. \quad (22)$$

(ii) For exponential cost the value function satisfies the following upper and lower bounds: On setting  $\tilde{\beta} = \beta e^{\kappa\alpha}$  and  $\tilde{x} = x - \tilde{\beta}$ , there are constants  $C_-, C_+$  such that, whenever  $x \geq \tilde{\beta}(e^2 + 1)$ ,

$$C_- + \frac{\kappa}{2e\beta} \frac{\tilde{x}^2}{\log(\tilde{x}) - (\kappa\alpha + 1)} \leq J^*(x) \leq C_+ + \frac{\kappa}{2e\beta} \tilde{x}^2$$

Part (i) of the above proposition exposes a connection between the fluid control policy and prior results about speed scaling obtained in the literature on worst-case algorithms, e.g. [5]. In particular, the optimal fluid control corresponds to a speed scaling scheme that is known to have a small competitive ratio.

Next, we can derive a lower bound on the difference  $c - c^\circ$  relatively easily.

**Lemma 2.**  $\mathcal{E}(x, u) \geq 0$  everywhere, giving  $c \geq c^\circ - \eta^\circ$ .

*Proof:* Convexity of  $J^*$  gives the bound,

$$J^*(Q(t+1)) - J^*(Q(t)) \geq \nabla J^*(Q(t)) \cdot (Q(t+1) - Q(t))$$

Consequently, for each  $x \in \mathbb{R}_+, u \in \mathbb{R}_+$  we have the lower bound,

$$\begin{aligned} P_u J^*(x) &= J^*(x) + \mathbb{E}_{x,u}[J^*(Q(1)) - J^*(Q(0))] \\ &\geq J^*(x) + \mathbb{E}_{x,u}[\nabla J^*(Q(0)) \cdot ((Q(1)) - Q(0))] \\ &= J^*(x) + \nabla J^*(x) \cdot (-u + \alpha) \end{aligned}$$

From the definition (20) this gives,

$$\mathcal{E}(x, u) \geq c(x, u) + \nabla J^*(x) \cdot (-u + \alpha)$$

Non-negativity follows from the TCOE (5).  $\square$

Further, we can derive an upper bound on  $c - c^\circ$  in two simple steps. We first write,

$$\underline{\mathcal{E}}(x) \leq \mathcal{E}(x, \phi^{F^*}(x)) \quad (23)$$

where  $\phi^{F^*}(x)$  is the optimal policy for the fluid model given in (7). Next we apply the second order Mean Value Theorem to bound  $\mathcal{E}$ . Given  $Q(0) = x$  and  $U(0) = u$  we have  $Q(1) = x - u + A(1)$ . For some random variable  $\bar{Q}$  between  $x$  and  $x - u + A(1)$  we have

$$\begin{aligned} \mathcal{D}_u J^*(x) &:= \mathbb{E}_{x,u}[J^*(Q(1)) - J^*(Q(0))] \\ &= \nabla J^*(x) \cdot (-u + \alpha) \\ &\quad + \frac{1}{2} \mathbb{E}[\nabla^2 J^*(\bar{Q}) \cdot (-u + A(1))^2] \end{aligned} \quad (24)$$

Proposition 1 states that the second derivative of  $J^*$  is non-increasing. Hence we can combine (24) with (23) to obtain,

$$\underline{\mathcal{E}}(x) \leq \frac{1}{2} \mathbb{E}[\nabla^2 J^*(x - \phi^{F^*}(x)) \cdot (-\phi^{F^*}(x) + A(1))^2]. \quad (25)$$

Lemma 3 provides an implication of this bound in the special case of quadratic cost.

**Lemma 3.** For polynomial cost (18) with  $\varrho = 2$ ,  $\beta = \frac{1}{2}$ , we have  $\underline{\mathcal{E}}(x) = \mathcal{O}(\sqrt{x})$ , and hence  $c(x, u) \leq c^\circ(x, u) + \mathcal{O}(\sqrt{x})$ .

*Proof:* The optimal policy is given in (22), giving  $\phi^{F^*}(x) = \mathcal{O}(\sqrt{x})$  in this special case. The formula (21)

gives  $\nabla^2 J^*(x) = \mathcal{O}(1/\sqrt{x})$ . The bound (25) then gives  $\underline{\mathcal{E}}(x) = \mathcal{O}(\sqrt{x})$ .  $\square$

Lemma 4 is an extension to the case of exponential cost. There is no space here for a proof.

**Lemma 4.** For exponential cost (18), with  $\beta = 1$ , we have  $\underline{\mathcal{E}}(x) \leq \kappa \log(x)^2$  for all  $x$  sufficiently large. For such  $x$  we have  $c(x, u) \leq c^\circ(x, u) - \eta^\circ + \kappa \log(x)^2$ .  $\square$

Hence, for quadratic or exponential cost, the fluid value function  $J^*$  can be interpreted as the relative value function for a cost function that approximates  $c(x, u)$ .

### B. The diffusion model

We next consider the diffusion model introduced in (9). We motivate the model using the second order Taylor series approximation (24). This continuous-time model will be used to obtain additional insight regarding the structure of  $h^*$ .

The ACOE for the diffusion model is similar to the total cost DP equation for the fluid model:

$$\min_{u \geq 0} \{c(x, u) + \mathcal{D}_u h^*(x)\} = \eta^* \quad (26)$$

where  $\eta^*$  is the average cost,  $h^*$  is called the relative value function, and  $\mathcal{D}_u$  denotes the usual differential generator. This is defined for  $C^2$  functions  $g: \mathbb{R}_+ \rightarrow \mathbb{R}_+$  via,

$$\mathcal{D}_u g(x) = \frac{d}{dx} g(x)(-u + \alpha) + \frac{1}{2} \sigma^2(u) \frac{d^2}{dx^2} g(x)$$

However, for a *reflected* diffusion the domain of the differential generator is restricted to those  $C^2$  functions satisfying the boundary condition,

$$\left. \frac{d}{dx} g(x) \right|_{x=0} = 0 \quad (27)$$

This is imposed so that the reflection term vanishes in the Ito formula:

$$dg(Q(t)) = f_g(Q(t), U(t)) dt + \sigma(U(t)) \frac{d}{dx} g(Q(t)) dN(t)$$

with  $f_g(x, u) = \mathcal{D}_u g(x)$ .

The variance term is selected so that the action of the differential generator on a smooth function will be similar to that of the discrete generator. The second order Taylor series expansion (24) suggests the value:

$$\sigma^2(u) = \mathbb{E}[(u - A(1))^2] = u^2 - 2\alpha u + m_A^2,$$

where  $m_A^2$  is the second moment of  $A(1)$ . We adopt this form in the remainder of this section.

Further, for the remainder of the section, we restrict to the case of quadratic cost:

$$c(x, u) = x + \frac{1}{2} u^2, \quad (28)$$

In this case the minimizer in (26) is given by,

$$\phi^*(x) := \frac{\nabla h^*(x) + \alpha \nabla^2 h^*(x)}{1 + \nabla^2 h^*(x)} \quad (29)$$

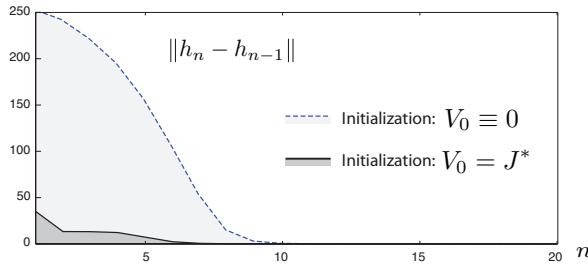


Fig. 2: The convergence of value iteration for the quadratic cost function (28). The error  $\|h_{n+1} - h_n\|$  converges to zero *much faster* when the algorithm is initialized using the fluid value function.

It can be shown that  $h^*$  is convex. Consequently, subject to the boundary condition (27), it follows that  $\phi^*(x) \geq 0$  for each  $x$ . Substituting (29) into (26) gives the fixed point equation,

$$x + \alpha \nabla h^* + \frac{1}{2} m_A^2 \nabla^2 h^* - \frac{(\alpha \nabla^2 h^* + \nabla h^*)^2}{2(1 + \nabla^2 h^*)} = \eta^*. \quad (30)$$

Although the cost function (28) does not satisfy  $c(0, \alpha) = 0$ , the TCOE (5) for the fluid model admits the solution,

$$J^*(x) = \alpha x + \frac{1}{3}[(2x + \alpha^2)^{3/2} - \alpha^3] \quad (31)$$

Furthermore, the function  $h^\circ(x) = J^*(x) + \frac{1}{2}x$  approximately solves the dynamic programming equation for the diffusion. In fact, it is straightforward to show that  $h^\circ(x)$  solves the ACOE for the diffusion exactly under a modified cost function:

$$c^\circ(x, u) = c(x, u) + \frac{1}{8} \left( \frac{y}{y+1} - 4 \frac{\sigma_A^2}{y} \right) + \eta^\circ,$$

where  $\sigma_A^2 = m_A^2 - \alpha^2$ , and  $y := (2x + \alpha^2)^{\frac{1}{2}}$ . The constant  $\eta^\circ$  is again arbitrary. Regardless of its value, the optimal average cost of  $c^\circ$  is equal to  $\eta^\circ$ . It is also easy to see that  $|c^\circ(x, u) - c(x, u)|$  is uniformly bounded over  $x$  and  $u$ .

The only issue that remains is the fact that  $h^\circ(x)$  does not satisfy the boundary condition (27) since

$$\nabla h^\circ(x) \Big|_{x=0} = 2\alpha + \frac{1}{2}.$$

This gap is resolved through an additional perturbation. Specifically, fix  $\vartheta > 0$ , and introduce the decaying exponential,

$$h^{\circ\circ} = h^\circ(x) - (2\alpha + \frac{1}{2})\vartheta e^{-x/\vartheta}$$

The gradient vanishes at the origin following this perturbation. This function solves the ACOE for the diffusion for a function  $c^{\circ\circ}$  which retains the property that  $c^{\circ\circ}(x, u) - c(x, u)$  is uniformly bounded.

Based on this form we are motivated to enlarge the basis to approximate the relative value function with  $\psi^1 = J^*$  and  $\psi^2(x) \equiv x$ .

## V. EXPERIMENTAL RESULTS

In this section we present results from experiments conducted for the speed scaling model described in Section III. Each of the value function approximations used in these experiments were based on insights obtained from the fluid and diffusion models.

In all of the numerical experiments described here the arrival process  $A$  is a scaled geometric distribution,

$$A(t) = \Delta_A G(t), \quad t \geq 1, \quad (32)$$

where  $\Delta_A > 0$  and  $G$  is geometrically distributed on  $\{0, 1, \dots\}$  with parameter  $p_A$ . The mean and variance of  $A(t)$  are given by, respectively,

$$m_A = \Delta_A \frac{p_A}{1 - p_A}, \quad \sigma_A^2 = \frac{p_A}{(1 - p_A)^2} \Delta_A^2. \quad (33)$$

### A. Value iteration

We begin by computing the actual solution to the average cost optimality equation using value iteration. This provides a reference for evaluating the proposed approach for TD learning. We restrict to the special case of the quadratic cost function given in (28) due to limited space. The arrival process is taken of the form (32), with  $p_A = 0.96$  and  $\Delta_A$  chosen so that the mean  $m_A$  is equal to unity:

$$1 = m_A = \Delta_A \frac{p_A}{1 - p_A} \quad \text{and} \quad \Delta_A = 1/24 \quad (34)$$

The state space is truncated for practical implementation of value iteration. In the experiments that follow we take  $X = \{\Delta_A m : m = 0, \dots, N_\ell\}$  with  $N_\ell = 480$ . The model becomes,

$$Q(t+1) = [Q(t) - U(t) + A(t+1)], \quad t \geq 0,$$

where  $[\cdot]$  represents projection to the interval  $[0, 20]$ , and  $U(t)$  is restricted to non-negative integer multiples of  $\Delta_A$ .

Let  $V_n$  denote the  $n$ th value function obtained. The approximate solution to the ACOE at stage  $n$  is taken to be the normalized value function  $h_n(x) = V_n(x) - V_n(0)$ ,  $x \in X$ . The convergence of  $\{h_n\}$  to  $h^*$  is illustrated in Fig. 2. The comparison of  $J^*$  and  $h^*$  shown in Fig. 1 was computed using this algorithm. Shown in Fig. 3 is the optimal policy and the  $(c, J^*)$ -myopic policy,  $\phi^J(x) = \arg \min_{0 \leq u \leq x} \{c(x, u) + P_u J^*(x)\}$ .

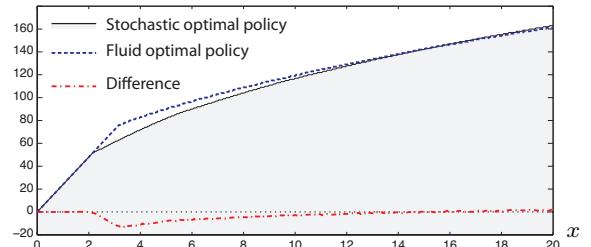


Fig. 3: The optimal policy compared to the  $(c, J^*)$ -myopic policy for the quadratic cost function (28).

### B. TD learning

We are now ready to apply TD learning to approximate the relative value function in the case of a specific policy. The policies considered here are taken to be the following translation of the optimal policy for the fluid model,

$$\phi_\diamond^{F*}(x) = \lfloor \min(x, \phi^F(x)) \rfloor, \quad x \in \mathbb{R}_+, \quad (35)$$

where here  $x$  is restricted to the lattice on which  $\mathbf{Q}$  evolves, and  $\lfloor a \rfloor$  indicates the nearest point on this lattice for  $a \in \mathbb{R}_+$ . In the next section we show how to combine TD learning and policy improvement in order to determine an approximately optimal solution.

We consider only polynomial costs due to space constraints. Additionally, we maintain the arrival distribution defined by (32), and the specification  $p_A = 0.96$  used in the previous subsection. We consider several values of  $\Delta_A$  to investigate the impact of variance on the estimation algorithm.

We take the following as the basis for TD learning

$$\psi_1(x) = J^*(x), \quad \psi_2(x) = x, \quad x \geq 0. \quad (36)$$

In the special case of quadratic cost, with  $c(x, u) = x + \frac{1}{2}u^2$ , this choice is motivated by the diffusion approximations presented in Sec. IV-B. We begin with results in this special case. Recall that the case of quadratic costs models the scenario of speed scaling in microprocessors.

The fluid value function  $J^*$  associated with the quadratic cost function (28) is given in (31). Fig. 1 shows a result obtained after 100,000 iterations of the LSTD algorithm. The initial condition was taken to be  $r(0) = (0, 0)^T$ . The value of the coefficient  $r_1^*$  corresponding to  $\psi_1 = J^*$  was found to be close to unity. Hence the approximate relative value function  $h^{r^*}$  is approximated by  $J^*$ , where  $r^*$  is the final value obtained from the LSTD algorithm. This conclusion is plainly illustrated in Fig. 1 where a plot of the function  $h^{r^*}$  is compared to the fluid value function  $J^*$  and the solution to the ACOE  $h^*$ .

We next turn to a different polynomial cost function. As a contrasting example we choose (18) with  $\varrho = 15$  and  $\beta = \varrho^{-1}$ . Shown on the left hand side of Fig. 4 are plots showing the evolution of the parameter estimates based on the LSTD algorithm with basis (36). The policy applied was the translation of the optimal policy for the fluid model given in (35). Shown on the right hand side is a comparison of the resulting approximation to Poisson's equation, and the fluid value function (21).

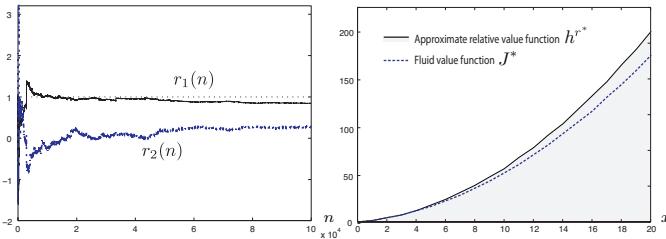


Fig. 4: The experiment illustrated in Fig. 1 was repeated for the cost function  $c(x, u) = x + [u - \alpha]_+^{15}/15$ .

### C. TD learning with policy improvement

So far, the TD learning algorithm was used to compute an approximation of the relative value function for the specific policy given in (35). In this section, we construct a policy using TD learning and policy improvement.

The policy iteration algorithm (PIA) is a method to construct an optimal policy through the following steps. The algorithm is initialized with a policy  $\phi^0$  and then the following operations are performed in the  $k$ th stage of the algorithm:

- (i) Given the policy  $\phi^k$ , find the solution  $h^k$  to Poisson's equation  $P_{\phi^k} h^k = h^k - c_k + \eta_k$ , where  $c_k(x) = c(x, \phi^k(x))$ , and  $\eta_k$  is the average cost.
- (ii) Update the policy via  $\phi^{k+1}(x) \in \arg \min_u \{c(x, u) + P_u h^k(x)\}$ .

In order to combine TD learning with PIA, the TDPIA algorithm considered replaces the first step with an application of the LSTD algorithm, resulting in an approximation  $h^{\text{TD}k}$  to the function  $h^k$ . The policy in (ii) is then taken to be  $\phi^{k+1}(x) \in \arg \min_u \{c(x, u) + P_u h^{\text{TD}k}(x)\}$ .

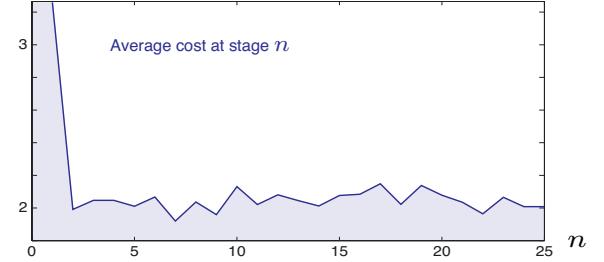


Fig. 5: Simulation result for TDPIA with the quadratic cost function (28), and basis  $\{\psi_1, \psi_2\} \equiv \{J^*, x\}$ .

We illustrate this approach in the case of the quadratic cost function (28), using the basis given in (36). The initial policy was taken to be  $\phi^0(x) = \min(x, 1)$ ,  $x \geq 0$ . Fig. 5 shows the estimated average cost in each of the twenty iterations of the algorithm. The algorithm results in a policy that is nearly optimal after just a few iterations.

## VI. CONCLUDING REMARKS

The main message of this paper is that idealized models (fluid and diffusion approximations) are useful choices when determining the function class for TD learning. This approach is applicable for control synthesis and performance approximation of Markov models in a wide range of applications. The motivation for this approach is a simple Taylor series argument that can be used to bound the difference between the relative value function  $h^*$  and the fluid value function  $J^*$ . Further, this approximation can be refined using a diffusion model.

To illustrate the application of this approach for TD learning, this paper focuses on a power management problem: dynamic speed scaling. This application reveals that this approach to approximation yields results that are remarkably accurate. In particular, numerical experiments revealed that (i) value iteration initialized using the fluid approximation results in much faster convergence, and (ii) policy iteration coupled with TD learning quickly converges to an approximately optimal policy when the fluid and diffusion models are considered in the construction of a basis. Further, the results studying the fluid model provided an interesting connection to worst-case analyses of the speed scaling problem.

Immediate extensions of this work to Q learning [18], [29] and to approximate dynamic programming based on linear programming [9], [10] are currently under investigation.

#### ACKNOWLEDGMENT

Financial support from the National Science Foundation (ECS-0523620 and CCF-0830511), ITMANET DARPA RK 2006-07284, and Microsoft Research is gratefully acknowledged.<sup>1</sup>

#### REFERENCES

- [1] IBM PowerPC.
- [2] Intel Xscale.
- [3] Susanne Albers and Hiroshi Fujiwara. Energy-efficient algorithms for flow time minimization. In *Lecture Notes in Computer Science (STACS)*, volume 3884, pages 621–633, 2006.
- [4] Nikhil Bansal, Tracy Kimbrel, and Kirk Pruhs. Speed scaling to manage energy and temperature. *J. ACM*, 54(1):1–39, March 2007.
- [5] Nikhil Bansal, Kirk Pruhs, and Cliff Stein. Speed scaling for weighted flow times. In *Proc. ACM-SIAM SODA*, pages 805–813, 2007.
- [6] D.P. Bertsekas and J. N. Tsitsiklis. *Neuro-Dynamic Programming*. Atena Scientific, Cambridge, Mass., 1996.
- [7] David P. Bunde. Power-aware scheduling for makespan and flow. In *Proc. ACM Symp. Parallel Alg. and Arch.*, 2006.
- [8] Ranveer Chandra, Ratul Mahajan, Thomas Moscibroda, Ramya Raghavendra, and Paramvir Bahl. A case for adapting channel width in wireless networks. In *Proc. ACM SIGCOMM*, Seattle, WA, August 2008.
- [9] D. P. de Farias and B. Van Roy. The linear programming approach to approximate dynamic programming. *Operations Res.*, 51(6):850–865, 2003.
- [10] D. P. Pucci de Farias and B. Van Roy. A cost-shaping linear program for average-cost approximate dynamic programming with performance guarantees. *Math. Oper. Res.*, 31(3):597–620, 2006.
- [11] Jennifer M. George and J. Michael Harrison. Dynamic control of a queue with adjustable service rate. *Operations Research*, 49(5):720–731, September 2001.
- [12] S. G. Henderson, S. P. Meyn, and V. B. Tadić. Performance evaluation and policy selection in multiclass networks. *Discrete Event Dynamic Systems: Theory and Applications*, 13(1-2):149–189, 2003. Special issue on learning, optimization and decision making (invited).
- [13] Sebastian Herbert and Diana Marculescu. Analysis of dynamic voltage/frequency scaling in chip-multiprocessors. In *Proc. ISLPED*, page 6, 2007.
- [14] A. Hoorfar and M. Hassani. Inequalities on the lambert  $w$  function and hyperpower function. *Journal of Inequalities in Pure and Applied Mathematics (JIPAM)*, 9(2), 2008.
- [15] Stefanos Kaxiras and Margaret Martonosi. *Computer Architecture Techniques for Power-Efficiency*. Morgan and Claypool, 2008.
- [16] S. Mannor, I. Menache, and N. Shimkin. Basis function adaptation in temporal difference reinforcement learning. *Annals of Oper. Res.*, 134(2):215–238, 2005.
- [17] Lykomidis Mastroleon, Daniel O'Neill, Benjamin Yolken, and Nick Bambos. Power aware management of packet switches. In *Proc. High-Perf. Interconn.*, 2007.
- [18] P. Mehta and S. Meyn. Machine learning and Pontryagin's Minimum Principle. Submitted to the 48th IEEE Conference on Decision and Control, December 16-18 2009.
- [19] S. P. Meyn. Stability and optimization of queueing networks and their fluid models. In *Mathematics of stochastic manufacturing systems (Williamsburg, VA, 1996)*, pages 175–199. Amer. Math. Soc., Providence, RI, 1997.
- [20] S. P. Meyn. Workload models for stochastic networks: Value functions and performance evaluation. *IEEE Trans. Automat. Control*, 50(8):1106–1122, August 2005.
- [21] S. P. Meyn. *Control Techniques for Complex Networks*. Cambridge University Press, Cambridge, 2007.
- [22] C.C. Moallemi, S. Kumar, and B. Van Roy. Approximate and data-driven dynamic programming for queueing networks. Submitted for publication., 2006.
- [23] S. Narendra et al. Ultra-low voltage circuits and processor in 180 nm to 90 nm technologies with a swapped-body biasing technique. In *Proc. IEEE Int. Solid-State Circuits Conf.*, page 8.4, 2004.
- [24] Kirk Pruhs, Patchrawat Uthaisombut, and Gerhard Woeginger. Getting the best response for your erg. In *Scandinavian Worksh. Alg. Theory*, 2004.
- [25] Kirk Pruhs, Rob van Stee, and Patchrawat Uthaisombut. Speed scaling of tasks with precedence constraints. In *Proc. Workshop on Approximation and Online Algorithms*, 2005.
- [26] R.S. Sutton and A.G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, on-line edition at <http://www.cs.ualberta.ca/~sutton/book/the-book.html> edition, 1998.
- [27] J. N. Tsitsiklis and B. Van Roy. An analysis of temporal-difference learning with function approximation. *IEEE Trans. Automat. Control*, 42(5):674–690, 1997.
- [28] M. H. Veatch. Approximate dynamic programming for networks: Fluid models and constraint reduction, 2004. Submitted for publication.
- [29] C. J. C. H. Watkins and P. Dayan. Q-learning. *Machine Learning*, 8(3-4):279–292, 1992.
- [30] W. Whitt. *Stochastic-process limits*. Springer Series in Operations Research. Springer-Verlag, New York, 2002.
- [31] A. Wierman, L. Andrew, and A. Tang. Power-aware speed scaling in processor sharing systems. In *Proc. of INFOCOM*, 2009.
- [32] L. Xie and P. R. Kumar. A network information theory for wireless communication: scaling laws and optimal operation. *IEEE Trans. on Info. Theory*, 50(5):748–767, 2004.
- [33] Francis Yao, Alan Demers, and Scott Shenker. A scheduling model for reduced CPU energy. In *Proc. IEEE Symp. Foundations of Computer Science (FOCS)*, pages 374–382, 1995.
- [34] Lin Yuan and Gang Qu. Analysis of energy reduction on dynamic voltage scaling-enabled systems. *IEEE Trans. Comput.-Aided Des. Integr. Circuits Syst.*, 24(12):1827–1837, December 2005.
- [35] Sushu Zhang and K. S. Catha. Approximation algorithm for the temperature-aware scheduling problem. In *Proc. IEEE Int. Conf. Comp. Aided Design*, pages 281–288, November 2007.
- [36] Yifan Zhu and Frank Mueller. Feedback EDF scheduling of real-time tasks exploiting dynamic voltage scaling. *Real Time Systems*, 31:33–63, December 2005.

<sup>1</sup>Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of NSF, DARPA, or Microsoft.