

# Gaussian processes for POMDP-based dialogue manager optimisation

Milica Gašić, *Member, IEEE*, Steve Young, *Fellow, IEEE*

**Abstract**—A partially observable Markov decision process (POMDP) has been proposed as a dialogue model that enables automatic optimisation of the dialogue policy and provides robustness to speech understanding errors. Various approximations allow such a model to be used for building real-world dialogue systems. However, they require a large number of dialogues to train the dialogue policy and hence they typically rely on the availability of a user simulator. They also require significant designer effort to hand-craft the policy representation. We investigate the use of Gaussian processes (GPs) in policy modelling to overcome these problems. We show that GP policy optimisation can be implemented for a real world POMDP dialogue manager, and in particular: 1) we examine different formulations of a GP policy to minimise variability in the learning process; 2) we find that the use of GP increases the learning rate by an order of magnitude thereby allowing learning by direct interaction with human users; and 3) we demonstrate that designer effort can be substantially reduced by basing the policy directly on the full belief space thereby avoiding ad hoc feature space modelling. Overall, the GP approach represents an important step forward towards fully automatic dialogue policy optimisation in real world systems.

**Index Terms**—Statistical dialogue systems, POMDP, Gaussian process

## I. INTRODUCTION

**S**POKEN dialogue systems enable human-computer interaction where the primary input is speech. As such they have innumerable benefits. However, building such systems to operate robustly is challenging as they are very sensitive to speech recognition errors and require designer effort to define their behaviour. The focus of this paper is automatic dialogue optimisation using a small amount of training data.

A partially observable Markov decision process (POMDP) has been proposed as a dialogue model which intrinsically deals with uncertainty from the recogniser, providing more robust system operation [1], [2], [3], [4], [5]. It assumes that the dialogue state  $s_t$  is only partially observable and depends on a noisy observation  $o_t$ . Since the dialogue state is unobservable, at every dialogue step a distribution over all states is maintained, which is called the *belief state*  $b_t$ . It takes values  $\mathbf{b} \in \mathcal{B}$ , where  $\mathcal{B}$  is a continuous space of dimensionality  $|\mathcal{S}|$ , namely  $[0, 1]^{|\mathcal{S}|}$ . The dialogue policy  $\pi$  then maps the belief state  $\mathbf{b}$  into an action  $a$  at every dialogue

turn,  $\pi(\mathbf{b}) = a$ . In each state, the system receives a reward  $r_t$ , which is a measure of how good the system action is for that particular state. The aim is then to find a policy such that the total accumulated reward at the end of the dialogue is maximised, which is measured by the  $Q$ -function.

The process of exact belief state updating becomes intractable for very large state spaces. However, there exist real-world dialogue systems based on the POMDP model which maintain an approximate belief state distribution in real-time throughout the dialogue. These include the Hidden Information State (HIS) system [6] and the Bayesian Update of Dialogue State (BUDS) system [5]. Policy optimisation, however, involves a large amount of training data and requires the designer to explicitly define the policy representation [7], [8]. This paper explores ways of overcoming these limitations.

## II. POMDP POLICY OPTIMISATION

Exact policy optimisation for POMDPs is intractable for everything but very simple cases [9]. Approximate POMDP solutions, as in the point-based value iteration algorithm [10], [11], are only suitable for relatively small state/action problems. However, a POMDP with discrete state and observation spaces can be viewed as an MDP with a continuous state space [9]. This allows standard MDP algorithms to be used for policy optimisation and this usually introduces the need for some form of parametrisation.

For example, it is proposed in [12] that the  $Q$ -function is parametrised as a linear combination of features from the belief state. Optimisation is then performed using a modified Sarsa algorithm. In the BUDS system, the policy is parametrised as a softmax function of features of the belief state [13]. Then, gradient methods and the natural-actor critic algorithm are used for policy optimisation. While this achieves tractability, in the case of the BUDS system it requires from  $10^5$  to  $10^6$  dialogues to train a policy, and this is only possible in interaction with a simulated user [8]. Moreover, the basis feature functions must be chosen by the designer and the solution is only optimal within the chosen basis. In [14], an algorithm is proposed that automatically selects useful features from an initial set of features predefined by a designer. An alternative approach is taken in [15] where Krylov iteration for lossless compression is used to compress the dialogue states, however the effectiveness of this approach in practice greatly depends on the definition of the reward function. A more recent thread of research has focused on parametric models of the  $Q$ -function where the uncertainty of the approximation can be encoded in the estimate of the parameters. Here,

Copyright (c) 2013 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

M. Gašić\* and S. Young are with the Department of Engineering, University of Cambridge, Trumpington St, Cambridge, CB2 1PZ, UK, e-mail: {mg436,sjy}@eng.cam.ac.uk

Manuscript received January 2013

a Gaussian distribution is placed over each parameter and Kalman filters are used for optimisation. It is shown that this leads to increased speed of training, however it still requires a set of basis functions [16].

In the HIS system [6], the belief space is heuristically mapped into a much smaller summary space, which is then discretised into a grid, allowing the Monte Carlo control algorithm to be used for the policy optimisation. This optimisation requires  $10^5$  dialogues, which is still too large for the policy to be optimised in direct interaction with humans and the forced use of a summary space again raises optimality questions.

Our aim is to reduce designer effort in defining a suitable policy representation and to speed up the process of policy optimisation to avoid reliance on user simulators. For this purpose, we propose non-parametric policy modelling which achieves operational efficiency by exploiting the similarities between different parts of the belief space. For instance, if the system action is reliably estimated in one belief state, that estimate should contribute to the unexplored parts of the belief space. In addition, it is not only important to estimate the system action for a given belief state, but also to provide a measure of how confident the system is about that estimate. This information is a useful indication of the extent to which different parts of the belief space have been explored during policy optimisation. A Gaussian process (GP) is proposed because it is a non-parametric model of Bayesian inference that has these desirable properties.

The next section introduces non-parametric policy modelling. It explains how the  $Q$ -function can be modelled as a Gaussian process and how a stochastic policy can be derived from this model that is well-suited for on-line learning. Then, in Section IV, experimental results using the BUDS dialogue manager are presented. Several issues concerning policy design are examined with the aim of reducing variability in the learning process. The resulting GP-based policy is shown to learn faster than a standard method and this is confirmed in a human evaluation using the Amazon MTurk crowdsourcing service [17]. A similar training scheme is conducted in interaction with humans via the Amazon MTurk service and is shown to perform significantly better than a simulator trained policy. The final set of experiments discuss policy training in the full belief space where it is shown that the GP approach will scale to handle the full space obviating the need for hand-crafted summary space mapping. Section V presents conclusions and directions for future work.

### III. NON-PARAMETRIC POLICY MODELLING

Non-parametric policy modelling avoids the limitations that occur when the solution is constrained by the chosen basis. This does not necessarily mean that it is parameter-free, but rather that the solution is not restricted by the choice of parameters. Instead, the choice of parameters only influences the speed at which the optimal solution is found.

A Gaussian process is a non-parametric Bayesian model used for function approximation. It has been successfully applied to reinforcement learning for continuous-space MDPs [18], [19], [20], [21]. An advantage of Bayesian approaches to policy optimisation is that they offer a principled

way of incorporating prior knowledge about the underlying task and this provides the potential to improve the learning rate. As already noted, it is important that the dependencies of different parts of the belief state space are taken into consideration during learning. Gaussian processes are able to incorporate the knowledge of these dependencies elegantly through the choice of a *kernel function*, the purpose of which is to describe correlations in different parts of the space. In addition, Gaussian processes can provide the uncertainty of the approximation by estimating the variance of the posterior.

In the GP approach, the size of the state space only impacts the evaluation of the kernel. Hence, it is not sensitive to the size of the POMDP summary space and indeed it can be used to model correlations directly in the full belief space. In this paper, the properties of the GP approach are evaluated in summary space in order to allow comparison with conventional gradient methods, but preliminary results using GP in the full belief space are also given.

This section provides an overview of the Gaussian process model of the  $Q$ -function based on the description given in [19]. Several policy formulations derived from the GP model are presented which are suitable for dialogue management. Issues of computational complexity and sparse approximation solutions are also discussed.

#### A. Gaussian process model for the $Q$ -function

This section describes how the  $Q$ -function can be modelled as a Gaussian process. As outlined in the introduction, a discrete-space POMDP can be viewed as a continuous-space MDP. In dialogue management, however, this space is often reduced to a summary space that contains both continuous and discrete variables [6]. Therefore, in the most general case, the approximation framework needs to support modelling of the  $Q$ -function in a multidimensional space that consists of both continuous and discrete variables. A Gaussian process allows such modelling. For now, however, an MDP with a full continuous belief state space  $\mathcal{B}$  is considered.

The discounted return  $R_t^\pi$  for time step  $t$  and a given policy  $\pi$  is the total accumulated reward acquired over time:

$$R_t^\pi = \sum_{i=0}^{\infty} \gamma^i r_{t+i+1}, \quad (1)$$

where  $r_t$  is the immediate reward at time step  $t$  and  $\gamma \in [0, 1]$ , is the discount factor. If the immediate reward is a random process, the discounted return is also a random process.

The discounted return can be written recursively as:

$$R_t^\pi = r_{t+1} + \gamma R_{t+1}^\pi. \quad (2)$$

The  $Q$ -function for policy  $\pi$  is the expectation of the discounted return given belief state  $\mathbf{b}$  and action  $a$  at time  $t$ , over all possible belief state sequences that can be generated:

$$Q^\pi(\mathbf{b}, a) = E_\pi(R_t | b(s_t) = \mathbf{b}, a_t = a). \quad (3)$$

Due to the stochasticity of transitions, the discounted return is a random variable and can be decomposed into a mean

$Q^\pi(\mathbf{b}, a)$  and a residual  $\Delta Q^\pi(\mathbf{b}, a)$ :

$$R_t^\pi(b(s_t) = \mathbf{b}, a_t = a) = Q^\pi(\mathbf{b}, a) + \Delta Q^\pi(\mathbf{b}, a). \quad (4)$$

Substituting  $R_t^\pi$  and  $R_{t+1}^\pi$  from (4) into (2) yields:

$$r_{t+1}(b(s_t) = \mathbf{b}, a_t = a) = Q^\pi(\mathbf{b}, a) - \gamma Q^\pi(\mathbf{b}', a') + \Delta Q^\pi(\mathbf{b}, a) - \gamma \Delta Q^\pi(\mathbf{b}', a'), \quad (5)$$

where  $b(s_{t+1}) = \mathbf{b}'$  is the next belief state and  $a' = \pi(\mathbf{b}')$  is the next action,  $a_{t+1} = a'$ .

Let  $\mathbf{B}_t = [(\mathbf{b}^0, a^0), \dots, (\mathbf{b}^t, a^t)]^\top$  be a sequence of  $t$  belief states and action samples<sup>1</sup> generated with policy  $\pi$ . Then, (5) becomes:

$$\begin{aligned} r^1 &= Q^\pi(\mathbf{b}^0, a^0) - \gamma Q^\pi(\mathbf{b}^1, a^1) \\ &\quad + \Delta Q^\pi(\mathbf{b}^0, a^0) - \gamma \Delta Q^\pi(\mathbf{b}^1, a^1) \\ r^2 &= Q^\pi(\mathbf{b}^1, a^1) - \gamma Q^\pi(\mathbf{b}^2, a^2) \\ &\quad + \Delta Q^\pi(\mathbf{b}^1, a^1) - \gamma \Delta Q^\pi(\mathbf{b}^2, a^2) \\ &\vdots \\ r^t &= Q^\pi(\mathbf{b}^{t-1}, a^{t-1}) - \gamma Q^\pi(\mathbf{b}^t, a^t) \\ &\quad + \Delta Q^\pi(\mathbf{b}^{t-1}, a^{t-1}) - \gamma \Delta Q^\pi(\mathbf{b}^t, a^t), \end{aligned} \quad (6)$$

where  $r^1, \dots, r^t$  are the acquired immediate rewards. This can be written in a more compact form as

$$\mathbf{r}_t = \mathbf{H}_t \mathbf{q}_t^\pi + \mathbf{H}_t \Delta \mathbf{q}_t^\pi, \quad (7a)$$

where

$$\mathbf{r}_t = [r^1, \dots, r^t]^\top \quad (7b)$$

$$\mathbf{q}_t^\pi = [Q^\pi(\mathbf{b}^0, a^0), \dots, Q^\pi(\mathbf{b}^t, a^t)]^\top, \quad (7c)$$

$$\Delta \mathbf{q}_t^\pi = [\Delta Q^\pi(\mathbf{b}^0, a^0), \dots, \Delta Q^\pi(\mathbf{b}^t, a^t)]^\top, \quad (7d)$$

$$\mathbf{H}_t = \begin{bmatrix} 1 & -\gamma & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \vdots & \ddots & \ddots & \vdots & \vdots \\ 0 & \dots & 0 & 1 & -\gamma \end{bmatrix}. \quad (7e)$$

If the transition probabilities and policy  $\pi$  do not change over time,  $Q^\pi(\mathbf{b}, a)$  is not a random variable, however, during the process of estimation  $\pi$  changes, so  $Q^\pi(\mathbf{b}, a)$  can be modelled as a random variable. By modelling  $Q^\pi(\mathbf{b}, a)$  as a Gaussian process,  $Q^\pi(\mathbf{b}, a) \sim \mathcal{GP}(0, k((\mathbf{b}, a), (\mathbf{b}, a)))$ , where the kernel  $k(\cdot, \cdot)$  is factored into separate kernels over the belief state and action spaces  $k_B(\mathbf{b}, \mathbf{b}')k_A(a, a')$  and  $\Delta Q^\pi(\mathbf{b}, a) \sim \mathcal{N}(0, \sigma^2)$  is Gaussian noise, Gaussian process regression [22] can be applied to find the posterior of  $Q^\pi(\mathbf{b}, a)$ , given a set of belief state-action pairs  $\mathbf{B}_t$  and the observed rewards  $\mathbf{r}_t$ :

$$\begin{aligned} Q^\pi(\mathbf{b}, a) | \mathbf{r}_t, \mathbf{B}_t &\sim \mathcal{N}(\bar{Q}(\mathbf{b}, a), \text{cov}((\mathbf{b}, a), (\mathbf{b}, a))), \\ \bar{Q}(\mathbf{b}, a) &= \mathbf{k}_t(\mathbf{b}, a)^\top \mathbf{H}_t^\top (\mathbf{H}_t \mathbf{K}_t \mathbf{H}_t^\top + \sigma^2 \mathbf{H}_t \mathbf{H}_t^\top)^{-1} \mathbf{r}_t, \\ \text{cov}((\mathbf{b}, a), (\mathbf{b}, a)) &= k((\mathbf{b}, a), (\mathbf{b}, a)) \\ &\quad - \mathbf{k}_t(\mathbf{b}, a)^\top \mathbf{H}_t^\top (\mathbf{H}_t \mathbf{K}_t \mathbf{H}_t^\top + \sigma^2 \mathbf{H}_t \mathbf{H}_t^\top)^{-1} \mathbf{H}_t \mathbf{k}_t(\mathbf{b}, a) \end{aligned} \quad (8)$$

where  $\mathbf{k}_t(\mathbf{b}, a) = [k((\mathbf{b}^0, a^0), (\mathbf{b}, a)), \dots, k((\mathbf{b}^t, a^t), (\mathbf{b}, a))]^\top$

<sup>1</sup>Random variables are denoted with a subscript, e.g.,  $a_t$  is a random variable at some time step  $t$ . Samples are denoted with a superscript e.g.,  $a^t$  is the action that was taken at time step  $t$ .

and  $\mathbf{K}_t = [\mathbf{k}_t((\mathbf{b}^0, a^0)), \dots, \mathbf{k}_t((\mathbf{b}^t, a^t))]$ .

This Gaussian process model of the  $Q$ -function exploits the relationship between *distributions* of the  $Q$ -function at differing time steps. This is in contrast to standard reinforcement learning algorithms which exploit the relationship between *values* of the  $Q$ -function at differing time steps. An important implication of this property is that prior knowledge about the  $Q$ -function can be incorporated in the prior distribution. Initially, the estimated  $Q$ -function distribution is just a zero-mean Gaussian process with the kernel function  $k((\mathbf{b}, a), (\mathbf{b}, a))$ , since no data has been observed. By time step  $t$ , the estimate is the posterior distribution given the set of observed rewards  $\mathbf{r}_t$  and associate belief state-action pairs  $\mathbf{B}_t$ .

It can be shown [23] that the marginal likelihood of the observed rewards is modelled by

$$\mathbf{r}_t | \mathbf{B}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{H}_t (\mathbf{K}_t + \sigma^2 \mathbf{I}) \mathbf{H}_t^\top). \quad (9)$$

If the kernel function is parametrised, this relationship can be used to optimise the kernel parameters as well as the noise parameter  $\sigma$  [24].

In the next section the derivation of a policy model from the Gaussian process model of the  $Q$ -function is discussed.

### B. Gaussian process-based policy model

The description so far has presented a Gaussian process model for the  $Q$ -function associated with a policy  $\pi$  given belief state-action pairs  $\mathbf{B}_t$  and rewards  $\mathbf{r}_t$  (8). What remains is to define the policy.

The simplest way of deriving such a policy to select the action associated with the highest mean of the  $Q$ -function:

$$\pi(\mathbf{b}) = \arg \max_a \bar{Q}(\mathbf{b}, a). \quad (10)$$

However, since the objective is to learn the  $Q$ -function associated with the optimal policy from interaction, the policy must exhibit some form of stochastic behaviour in order to explore alternatives during the process of learning. One way of defining a policy that can be optimised on-line is to use an  $\epsilon$ -greedy approach. This requires setting an additional parameter  $\epsilon$  which balances how often an action is taken according to the current best estimate of the  $Q$ -function mean (the *exploitation* stage) compared to how often an action is taken randomly (the *exploration* stage). This policy is defined as:

$$\pi(\mathbf{b}) = \begin{cases} \arg \max_a \bar{Q}(\mathbf{b}, a) & \text{with prob } 1 - \epsilon, \\ \text{random } a & \text{with prob } \epsilon. \end{cases} \quad (11)$$

However, such random exploration can be inefficient since not all parts of belief-action space are equally informative. Using active learning to provide more efficient exploration should yield faster learning [25]. Instead of selecting actions randomly, active learning selects actions according to some *utility function* which normally includes some measure of information gain through which various heuristics can be incorporated [26]. Gaussian process reinforcement learning provides a measure of uncertainty at each point in belief-action space. This uncertainty can then be used directly in the active learning utility function [21]. This enables the model

to explore the parts of the space it is less certain about. Therefore, during exploration, actions are chosen according to the variance of the GP estimate for the  $Q$ -function and during exploitation, actions are chosen according to the mean:

$$\pi(\mathbf{b}) = \begin{cases} \arg \max_a \bar{Q}(\mathbf{b}, a) & \text{with prob } 1 - \epsilon, \\ \arg \max_a \text{cov}((\mathbf{b}, a), (\mathbf{b}, a)) & \text{with prob } \epsilon. \end{cases} \quad (12)$$

Both of the policies defined in (11) and (12) require the proportions of exploration and exploitation to be balanced manually by tuning the parameter  $\epsilon$ . Since the Gaussian process for the  $Q$ -function defines a Gaussian distribution for every belief state-action pair (8), when a new belief point  $\mathbf{b}$  is encountered, for each action  $a \in \mathcal{A}$ , there is a Gaussian distribution  $\hat{Q}(\mathbf{b}, a) \sim \mathcal{N}(\bar{Q}(\mathbf{b}, a), \text{cov}((\mathbf{b}, a), (\mathbf{b}, a)))$ . Sampling from these Gaussian distributions gives  $Q$ -values from which the action with the highest sampled  $Q$ -value can be selected:

$$\pi(\mathbf{b}) = \arg \max_a \left\{ \hat{Q}(\mathbf{b}, a) : a \in \mathcal{A} \right\}. \quad (13)$$

If two actions have overlapping distributions then sampling from these distributions will result in the system randomly switching between them. Thus, in contrast to the  $\epsilon$ -greedy policies defined by (11) and (12), this approach maps the GP approximation of the  $Q$ -function into a stochastic policy model. Note however that all three policies are *greedy* in the sense that as either  $\epsilon$  or the variance tends to zero, they tend to select the locally optimal action. The performance of different policy models is discussed further below in Section IV-B.

### C. Gaussian process-based policy optimisation

If the  $Q$ -function of a greedy policy is estimated on-line using temporal differences, then it will in the limit tend to the  $Q$ -function of the optimal policy. A commonly used embodiment of this idea for MDPs is the Sarsa algorithm which iteratively updates the  $Q$ -function on-line using the rule:

$$Q(\mathbf{b}, a) \leftarrow Q(\mathbf{b}, a) + \alpha[r_{t+1}(b_t = \mathbf{b}, a_t = a) + \gamma Q(\mathbf{b}', a') - Q(\mathbf{b}, a)]. \quad (14)$$

If the  $Q$ -function correction provided by the right hand side of (14) is replaced by the update of the GP-based  $Q$ -function posterior given by (8), the GP-Sarsa algorithm is obtained. This is shown in Fig. 1 amended for episodic reinforcement learning and hence applicable to dialogue optimisation.

In practice, however, the exact computation of the posterior using (8) is intractable as it requires the inversion of a matrix of dimensionality  $t$ . A solution is discussed in the next section.

### D. Computational complexity and sparse approximation

Due to the matrix inversion in (8), the computational complexity of calculating the  $Q$ -function posterior is  $O(t^3)$ , where  $t$  is the number of data points. In the case of a dialogue system, this will be equal to the total number of turns, summed over all dialogues, which poses a serious computational problem.

One solution to this problem is to restrict the set of data points used for the  $Q$ -function approximation. This has the obvious drawback that useful information will be discarded

Fig. 1. Episodic GP-Sarsa (without computational constraints)

```

1: for each episode do
2:   Initialise  $\mathbf{b}$ 
3:   if first episode then
4:     Choose  $a$  arbitrary
5:      $\mathbf{B}_0 \leftarrow (\mathbf{b}, a)$ 
6:      $\mathbf{K}_0 \leftarrow k((\mathbf{b}, a), (\mathbf{b}, a))$ 
7:      $\mathbf{H}_0 \leftarrow [1 \quad -\gamma]$ 
8:   else
9:     if initial step then
10:       Choose  $a \leftarrow \pi(\mathbf{b})$  (11), (12) or (13)
11:     end if
12:   end if
13:   for each step in the episode do
14:     Take  $a$ , observe  $r'$ , update  $\mathbf{b}'$ 
15:     if non-terminal step then
16:       Choose new action  $a' \leftarrow \pi(\mathbf{b}')$  (11), (12) or (13)
17:        $\mathbf{B}_{t+1} \leftarrow \begin{bmatrix} \mathbf{B}_t & (\mathbf{b}', a') \end{bmatrix}$ 
18:        $\mathbf{H}_{t+1} \leftarrow \begin{bmatrix} \mathbf{H}_t & \mathbf{0} \\ \mathbf{u}^\top & -\gamma \end{bmatrix}$ , where  $\mathbf{u} = [0 \quad 1]^\top$ 
19:        $\mathbf{K}_{t+1} \leftarrow \begin{bmatrix} \mathbf{K}_t & \mathbf{k}_t(\mathbf{b}', a') \\ \mathbf{k}_t(\mathbf{b}', a') & k(\mathbf{b}', a', \mathbf{b}', a') \end{bmatrix}$ 
20:     else
21:        $\mathbf{B}_{t+1} \leftarrow \mathbf{B}_t$ ,  $\mathbf{K}_{t+1} \leftarrow \mathbf{K}_t$ 
22:        $\mathbf{H}_{t+1} \leftarrow \begin{bmatrix} \mathbf{H}_t \\ \mathbf{u}^\top \end{bmatrix}$ , where  $\mathbf{u} = [0 \quad 1]^\top$ 
23:     end if
24:      $\mathbf{r}_{t+1} \leftarrow [r_t \quad r']$ 
25:     Update  $Q$ -function posterior  $Q^\pi | \mathbf{r}_{t+1}, \mathbf{B}_{t+1}$  (8)
26:     if non-terminal step then
27:        $\mathbf{b} \leftarrow \mathbf{b}'$ ,  $a \leftarrow a'$ 
28:     end if
29:   end for
30: end for

```

and this is exacerbated when the space is continuous since the exact same state is never visited twice. Sparse approximation methods for Gaussian processes aim to take into account all data points whilst reducing the computational complexity.

A significant research effort has been invested into solving this problem, and a number of methods have been developed [27]. However, these rely on the pre-selection of a fixed set of  $m$  support points. This reduces the complexity of calculating the posterior from  $O(t^3)$  to  $O(tm^2)$  and if the number of support points is significantly smaller than the number of data points, this approach is very effective for reducing the computational cost. This class of methods is not useful for learning on-line, since the support points cannot be determined *a priori*. An alternative algorithm which approximates the Gaussian process without first obtaining a set of support points is the kernel span sparsification method described in [28]. In this case, a representative set of data points is acquired as the belief-action space is traversed during interaction with the user. It is assumed that the environment dynamics do not change over time hence obviating the need to remove support points and thereby avoiding the difficulties discussed in [29].

A kernel function can be thought of as the dot product of a (potentially infinite) set of feature functions  $k((\mathbf{b}, a), (\mathbf{b}, a)) = \langle \phi(\mathbf{b}, a), \phi(\mathbf{b}, a) \rangle$  where  $\phi(\mathbf{b}, a)$  is the vector of feature functions  $[\phi_1(\mathbf{b}, a), \phi_2(\mathbf{b}, a), \dots]^\top$ . Any linear combination of feature vectors  $\{\phi(\mathbf{b}^0, a^0), \dots, \phi(\mathbf{b}^t, a^t)\}$  for a given set of points  $\{(\mathbf{b}^0, a^0), \dots, (\mathbf{b}^t, a^t)\}$  is called the *kernel span*. The aim then is to find the subset of points that approximates this kernel span. These points are called the *representative points* and the set of representative points is called the *dictionary*  $\mathcal{D} = \{(\tilde{\mathbf{b}}^0, \tilde{a}^0), \dots, (\tilde{\mathbf{b}}^m, \tilde{a}^m)\}$ .

A sparsification parameter  $\nu$  places a threshold on the squared distance between the feature function span at the representative points, and the true feature function value at each visited point:

$$\min_{\mathbf{g}_t} \left\| \sum_{j=0}^m g_{tj} \phi(\tilde{\mathbf{b}}^j, \tilde{a}^j) - \phi(\mathbf{b}^t, a^t) \right\|^2 \leq \nu, \quad (15)$$

where  $(\mathbf{b}^t, a^t)$  is the current point,  $\mathbf{g}_t = [g_{t1}, \dots, g_{tm}]$  is a vector of coefficients and  $m$  is the size of the current dictionary,  $\mathcal{D} = \{(\tilde{\mathbf{b}}^0, \tilde{a}^0), \dots, (\tilde{\mathbf{b}}^m, \tilde{a}^m)\}$ . This is equivalent to [28]:

$$\min_{\mathbf{g}_t} \left( k((\mathbf{b}^t, a^t), (\mathbf{b}^t, a^t)) - \tilde{\mathbf{K}}_{t-1}(\mathbf{b}^t, a^t)^\top \mathbf{g}_t \right) \leq \nu, \quad (16)$$

where  $\tilde{\mathbf{K}}_{t-1}(\mathbf{b}^t, a^t) = [k((\mathbf{b}^t, a^t), (\tilde{\mathbf{b}}^0, \tilde{a}^0)), \dots, k((\mathbf{b}^t, a^t), (\tilde{\mathbf{b}}^m, \tilde{a}^m))]^\top$ . It can also be shown that the expression on the left side of (16) is minimised when  $\mathbf{g}_t = \tilde{\mathbf{K}}_{t-1}^{-1} \tilde{\mathbf{K}}_{t-1}(\mathbf{b}^t, a^t)$ , where  $\tilde{\mathbf{K}}_{t-1}$  is the Gram matrix of the current set of representative points. If the threshold  $\nu$  is exceeded then  $(\mathbf{b}^t, a^t)$  is added to the dictionary, otherwise the dictionary stays the same. This constitutes the sparsification criterion and it allows the posterior to be approximated with bounded complexity as follows.

Since the kernel function is the dot product of the feature functions, the Gram matrix is  $\mathbf{K}_t = \Phi_t^\top \Phi_t$  where  $\Phi_t = [\phi(\mathbf{b}^0, a^0), \dots, \phi(\mathbf{b}^t, a^t)]$ . The feature function values for each point are then approximated as the linear combination of the representative points  $\phi(\mathbf{b}^i, a^i) \approx \sum_{j=1}^m g_{ij} \phi(\tilde{\mathbf{b}}^j, \tilde{a}^j)$ , for all  $i \in 0, \dots, t$ . Also, the Gram matrix is approximated as

$$\mathbf{K}_t = \Phi_t^\top \Phi_t \approx \mathbf{G}_t \tilde{\mathbf{K}}_t \mathbf{G}_t^\top, \quad (17)$$

where  $\mathbf{G}_t = [\mathbf{g}_1, \dots, \mathbf{g}_t]$ . Similarly,  $\mathbf{k}_t(\mathbf{b}, a) \approx \mathbf{G}_t \tilde{\mathbf{k}}_t(\mathbf{b}, a)$ .

This allows the posterior (8) to be approximated as:

$$\begin{aligned} Q(\mathbf{b}, a) | \mathbf{B}_t, \mathbf{r}_t &\sim \mathcal{N} \left( \tilde{Q}(\mathbf{b}, a), \tilde{\text{cov}}((\mathbf{b}, a), (\mathbf{b}, a)) \right), \\ \tilde{Q}(\mathbf{b}, a) &= \tilde{\mathbf{k}}_t(\mathbf{b}, a)^\top (\tilde{\mathbf{H}}_t^\top (\tilde{\mathbf{H}}_t \tilde{\mathbf{K}}_t \tilde{\mathbf{H}}_t^\top + \sigma^2 \tilde{\mathbf{H}}_t \tilde{\mathbf{H}}_t^\top)^{-1} \mathbf{r}_t), \\ \tilde{\text{cov}}((\mathbf{b}, a), (\mathbf{b}, a)) &= k((\mathbf{b}, a), (\mathbf{b}, a)) \\ &\quad - \tilde{\mathbf{k}}_t(\mathbf{b}, a)^\top (\tilde{\mathbf{H}}_t^\top (\tilde{\mathbf{H}}_t \tilde{\mathbf{K}}_t \tilde{\mathbf{H}}_t^\top + \sigma^2 \tilde{\mathbf{H}}_t \tilde{\mathbf{H}}_t^\top)^{-1} \tilde{\mathbf{H}}_t) \tilde{\mathbf{k}}_t(\mathbf{b}, a), \end{aligned} \quad (18)$$

where  $\tilde{\mathbf{H}}_t = \mathbf{H}_t \mathbf{G}_t$ . It can be shown that this reduces the complexity to  $O(tm^2)$ , where  $m$  is the size of the dictionary, see [19] for details.

This sparsification method has a drawback that it turns a non-parametric method into a parametric method. More precisely, the kernel is now approximated using only a limited number of points, which is equivalent to defining a functional

basis for the kernel function. This can be shown to be equivalent to parametrising the  $Q$ -function [19]. This may limit the accuracy of the solution, however, the fact that the basis is chosen dynamically, without an initial set of basis functions, still allows for more appropriate functions than when using a fixed basis defined by a designer.

An advantage of this sparsification approach is that it enables non-positive definite kernel functions to be used in the approximation, for example see [24]. This is due to the fact that the sparsification method essentially changes the kernel function in a manner which ensures that the approximated Gram matrix is positive definite and this is sufficient to guarantee that the model remains well-defined [30].

This sparsification method allows observations to be processed sequentially, in direct interaction with the user, and it can be incorporated directly into the GP-Sarsa algorithm out-lined in Fig. 1. The full GP-Sarsa algorithm with kernel sparsification is given in Fig. 2.

## IV. EXPERIMENTS

The previous sections have introduced the theory of Gaussian process-based POMDP policy optimisation. Here we examine their application in practice and investigate three important research issues. Firstly, section IV-B addresses the detailed design trade-offs involved in practical implementation of a GP-based dialogue policy. Secondly, section IV-C investigates the speed of the GP policy optimisation process and whether in practice it can be performed in direct interaction with real users. Finally, section IV-D addresses the extent to which design time effort can be reduced by applying a GP-policy directly to the full belief space.

### A. Experimental set-up

A spoken dialogue system providing restaurant information for Cambridge UK was used as a test-bed for all of the reported experiments. This system uses the BUDS dialogue manager in which beliefs are modelled by a Bayesian network. Training and testing was provided by an agenda-based simulated user and by real users recruited using the Amazon MTurk service.

#### 1) Bayesian Update of Dialogue State dialogue manager:

The Bayesian Update of Dialogue State (BUDS) dialogue manager is a POMDP-based dialogue manager [13] whose belief state consists of the marginal posterior probability distribution over hidden nodes in a Bayesian network. Each concept in the dialogue, eg. *area*, *food-type*, *address* is represented by a pair of hidden nodes recording the history and the goal. The history nodes define possible states eg. *system-informed*, *user-informed* and the goal nodes define possible values for a particular concept, eg. *Chinese*, *Indian*.

In order to train an optimal policy the belief space  $\mathcal{B}$  is mapped into a summary space  $\mathcal{C}$ . Likewise, the action space is mapped into a smaller scale summary action space  $\mathcal{A}$ . This summary space  $\mathcal{C} \times \mathcal{A}$  is then mapped into a feature space  $\mathcal{F}$ . This space is then used to produce a parametric representation of the policy. A weighted set of these basis functions can then

Fig. 2. Episodic GP-Sarsa

```

1: Initialise  $\tilde{\boldsymbol{\mu}} \leftarrow \mathbf{0}$ ,  $\tilde{\mathbf{C}} \leftarrow \mathbf{0}$ ,  $\tilde{\mathbf{c}} \leftarrow \mathbf{0}$ ,  $d \leftarrow 0$ ,  $1/v \leftarrow 0$ 
2: for each episode do
3:   Initialise  $\mathbf{b}$ 
4:   if first episode then
5:     Choose  $a$  arbitrary,  $\mathcal{D} = \{(\mathbf{b}, a)\}$ ,  $\tilde{\mathbf{K}}^{-1} \leftarrow 1/k((\mathbf{b}, a), (\mathbf{b}, a))$ 
6:   else
7:     if initial step then
8:       Choose  $a \leftarrow \pi(\mathbf{b})$ 
9:     end if
10:  end if
11:   $\tilde{\mathbf{c}} \leftarrow \mathbf{0}$  {size  $|\mathcal{D}|$ },  $d \leftarrow 0$ ,  $1/v \leftarrow 0$ 
12:   $\mathbf{g} \leftarrow \tilde{\mathbf{K}}^{-1} \mathbf{k}(\mathbf{b}, a)$ ,  $\delta \leftarrow k((\mathbf{b}, a), (\mathbf{b}, a)) - \tilde{\mathbf{k}}(\mathbf{b}, a)^\top \mathbf{g}$ 
13:  if  $\delta > \nu$  then
14:     $\mathcal{D} \leftarrow \{(\mathbf{b}, a)\} \cup \mathcal{D}$ ,  $\tilde{\mathbf{K}}^{-1} \leftarrow \frac{1}{\delta} \begin{bmatrix} \delta \tilde{\mathbf{K}}^{-1} + \mathbf{g} \mathbf{g}^\top & -\mathbf{g} \\ -\mathbf{g}^\top & 1 \end{bmatrix}$ 
15:     $\mathbf{g} \leftarrow [0, \dots, 0, 1]^\top$  {size  $|\mathcal{D}|$ },  $\tilde{\boldsymbol{\mu}} \leftarrow \begin{bmatrix} \tilde{\boldsymbol{\mu}} \\ 0 \end{bmatrix}$ ,  $\tilde{\mathbf{C}} \leftarrow \begin{bmatrix} \tilde{\mathbf{C}} & \mathbf{0} \\ \mathbf{0}^\top & 0 \end{bmatrix}$ ,  $\tilde{\mathbf{c}} \leftarrow \begin{bmatrix} \tilde{\mathbf{c}} \\ 0 \end{bmatrix}$ 
16:  end if
17:  for each step in the episode do
18:    Take  $a$ , observe  $r'$ , update  $\mathbf{b}'$ 
19:    if non-terminal step then
20:      Choose new action  $a' \leftarrow \pi(\mathbf{b}')$ 
21:       $\mathbf{g}' \leftarrow \tilde{\mathbf{K}}^{-1} \mathbf{k}(\mathbf{b}', a')$ ,  $\delta \leftarrow k((\mathbf{b}', a'), (\mathbf{b}', a')) - \mathbf{k}(\mathbf{b}', a')^\top \mathbf{g}'$ 
22:       $\Delta \mathbf{k} \leftarrow \mathbf{k}(\mathbf{b}, a) - \gamma \tilde{\mathbf{k}}(\mathbf{b}', a')$ 
23:    else
24:       $\mathbf{g}' \leftarrow [0]$  {size  $|\mathcal{D}|$ },  $\delta \leftarrow 0$ ,  $\Delta \tilde{\mathbf{k}} \leftarrow \tilde{\mathbf{k}}(\mathbf{b}, a)$ 
25:    end if
26:     $d \leftarrow \frac{\gamma \sigma^2}{v} d + r' - \Delta \tilde{\mathbf{k}}^\top \tilde{\boldsymbol{\mu}}$ 
27:    if  $\delta > \nu$  and non-terminal step then
28:       $\mathcal{D} \leftarrow \{(\mathbf{b}', a')\} \cup \mathcal{D}$ ,  $\tilde{\mathbf{K}}^{-1} \leftarrow \frac{1}{\delta} \begin{bmatrix} \delta \tilde{\mathbf{K}}^{-1} + \mathbf{g} \mathbf{g}^\top & -\mathbf{g} \\ -\mathbf{g}^\top & 1 \end{bmatrix}$ 
29:       $\mathbf{g}' \leftarrow [0, \dots, 0, 1]^\top$ ,  $\tilde{\mathbf{h}} \leftarrow [\mathbf{g}^\top, -\gamma]^\top$ 
30:       $\Delta k_{tt} \leftarrow \mathbf{g}'^\top (\mathbf{k}(\mathbf{b}, a) - 2\gamma \tilde{\mathbf{k}}(\mathbf{b}', a')) + \gamma^2 k((\mathbf{b}', a'), (\mathbf{b}', a'))$ 
31:       $\tilde{\mathbf{c}}' \leftarrow \frac{\gamma \sigma^2}{v} \begin{bmatrix} \tilde{\mathbf{c}} \\ 0 \end{bmatrix} + \tilde{\mathbf{h}} - \begin{bmatrix} \tilde{\mathbf{C}} \Delta \tilde{\mathbf{k}} \\ 0 \end{bmatrix}$ 
32:       $v \leftarrow (1 + \gamma^2) \sigma^2 + \Delta k_{tt} - \Delta \tilde{\mathbf{k}}^\top \tilde{\mathbf{C}} \Delta \tilde{\mathbf{k}} + \frac{2\gamma \sigma^2}{v} \tilde{\mathbf{c}} \Delta \tilde{\mathbf{k}} - \frac{\gamma^2 \sigma^4}{v}$ 
33:       $\tilde{\boldsymbol{\mu}} \leftarrow \begin{bmatrix} \tilde{\boldsymbol{\mu}} \\ 0 \end{bmatrix}$ ,  $\tilde{\mathbf{C}} \leftarrow \begin{bmatrix} \tilde{\mathbf{C}} & \mathbf{0} \\ \mathbf{0}^\top & 0 \end{bmatrix}$ 
34:    else
35:       $\tilde{\mathbf{h}} \leftarrow \mathbf{g} - \gamma \mathbf{g}'$ ,  $\tilde{\mathbf{c}}' \leftarrow \frac{\gamma \sigma^2}{v} \tilde{\mathbf{c}} + \tilde{\mathbf{h}} - \tilde{\mathbf{C}} \Delta \tilde{\mathbf{k}}$ 
36:      if non-terminal step then
37:         $v \leftarrow (1 + \gamma^2) \sigma^2 + \Delta \tilde{\mathbf{k}}^\top (\tilde{\mathbf{c}}' + \frac{\gamma \sigma^2}{v} \tilde{\mathbf{c}}) - \frac{\gamma^2 \sigma^4}{v}$ 
38:      else
39:         $v \leftarrow \sigma^2 + \Delta \tilde{\mathbf{k}}^\top (\tilde{\mathbf{c}}' + \frac{\gamma \sigma^2}{v} \tilde{\mathbf{c}}) - \frac{\gamma^2 \sigma^4}{v}$ 
40:      end if
41:    end if
42:     $\tilde{\boldsymbol{\mu}} \leftarrow \tilde{\boldsymbol{\mu}} + \frac{\delta}{v} d$ ,  $\tilde{\mathbf{C}} \leftarrow \tilde{\mathbf{C}} + \frac{1}{v} \tilde{\mathbf{c}}' \tilde{\mathbf{c}}'^\top$ ,  $\tilde{\mathbf{c}} \leftarrow \tilde{\mathbf{c}}'$ ,  $\mathbf{g} \leftarrow \mathbf{g}'$ 
43:    if non-terminal step then
44:       $\mathbf{b} \leftarrow \mathbf{b}'$ ,  $a \leftarrow a'$ 
45:    end if
46:  end for
47: end for

```

be cast into a probability via a soft-max function

$$\pi(a|\mathbf{c}; \theta) = \frac{e^{\theta \cdot f_a(\mathbf{c})}}{\sum_a e^{\theta \cdot f_a(\mathbf{c})}}. \quad (19)$$

where  $f_a(\mathbf{c})$  is the vector of features derived from summary state  $\mathbf{c}$  for action  $a$ . The policy is then optimised using the natural actor critic (NAC) algorithm which is a gradient-based method [31]. In operation, system responses are generated by sampling (19) and then heuristically mapping the summary action back into a full system action. Once the system has chosen the summary action, the mapping to a full action is usually straightforward. For instance, if the action is `ConfirmArea` the system simply confirms the most likely value of the *area* slot.

Mappings  $\mathcal{B} \rightarrow \mathcal{C}$  and  $\mathcal{C} \times \mathcal{A} \rightarrow \mathcal{F}$ , however, require a significant amount of hand-crafting. Also, due to the parametric policy representation, the solution is only optimal within the given basis. Finally, gradient-based optimisation methods are inherently slow which prevents direct on-line optimisation with real users. Here we show that the GP-Sarsa algorithm can overcome these limitations.

2) *The Cambridge restaurant domain*: The Cambridge restaurant domain consists of a selection of approximately 150 restaurants, referred to as *entities*, which have been automatically extracted from various web-based sources. Each restaurant has 8 attributes (slots) and this results in a belief space consisting of 25 concepts where each concept takes from 3 to 150 values and each value has a probability in  $[0, 1]$ .<sup>2</sup> The summary space is formed from discrete, as well as continuous, features. Continuous features represent the entropy of the distribution for each hidden node in the Bayesian network. In addition, there are a few discrete features that correspond to history nodes, e.g. if the highest probability method that the user used to inform the system was *by name of the venue* or if the highest probability discourse act used was *repeat*. Also, there is a count of the number of top probability goal values that are greater than 0.8. Finally, there is a list which defines the order in which the goal slots can be accepted. In total, there are 129 features. The summary action space consists of 16 summary actions.

3) *The agenda-based simulated user*: In training and testing, an agenda-based user simulator was used [32], [33]. The user state is factored into an *agenda* and a *goal*. The goal ensures that the user simulator exhibits consistent, goal-directed behaviour. The role of the agenda is to elicit the dialogue acts that are needed for the user simulator to fulfil the goal. Both the goal and the agenda are dynamically updated throughout the dialogue. These updates are controlled by deterministic or stochastic *decision points*, allowing a wide spread of realistic dialogues to be generated.

In addition, an error model was used to add confusions to the simulated user input such that it resembles those found in real data [34]. The output of the error model is an N-best list of possible user responses. The length of this list was set

<sup>2</sup>It should be emphasised that the system is not sensitive to the number of values in a slot or the number of entities in the database since only values and entities which are actually mentioned in the dialogue are modelled at run time.

to 10 and the confusion rate was set to 15%, which means that 15% of time the true hypothesis is not in the N-best list. Intermediate experiments showed that these settings match the confusions typically found in real data.

The reward function was set to give a reward of 20 for successful dialogues, zero otherwise, less the number of dialogue turns. A successful dialogue is one in which the user goal is fulfilled. If the user goal changes, success depends on satisfying the final the user goal. The discount factor  $\gamma$  is set to 1 and the dialogue length is limited to 30 turns.

4) *The evaluation schedule:* When evaluating a policy optimisation technique one is typically interested in the performance of the resulting policy. However, it is also important to know how many dialogue sessions are needed to reach that performance. Policy optimisation is a random process itself. A lucky choice of actions at the beginning may lead to very good performance, whereas an unlucky choice of actions can slow down the process of learning. Therefore, it is important to examine the variability in performance during training.

For this reason, in all the user simulator-based experiments reported here, unless stated otherwise, 10 training sessions were seeded with different random seeds and after every 1000 dialogues the performance of the partially trained policies were evaluated. Each evaluation was performed with 1000 dialogues ensuring that the simulated user had a different random seed to the one used in training. The results were then averaged and presented with one standard error.

5) *The MTurk Amazon Service:* In order to evaluate different policies with real people we recruited subjects using the Amazon MTurk crowd-sourcing service in a set-up similar to [35]. The BUDS dialogue manager was incorporated in a live telephone-based spoken dialogue system in which human users were assigned specific tasks which involved finding restaurants that have particular features. To elicit more complex dialogues, users were sometimes asked to find more than one restaurant, or a restaurant that does not exist. In the latter case, they were asked to change one of their constraints, for example find a Chinese restaurant instead of a Vietnamese one. After each dialogue, users filled in a feedback form indicating whether they thought the dialogue was successful. This rating was then used to calculate the average success and reward.

## B. Design of Gaussian process-based policy optimisation

To use the GP-Sarsa algorithm, the kernel function, the noise parameter  $\sigma$  and the sparsification parameter  $\nu$  need to be defined. Whilst the choice of the kernel function and the choice of  $\sigma$  do not greatly affect the resulting policy, they do influence the speed of learning. In addition, the parameter  $\nu$  controls the approximation of the kernel function, essentially turning a non-parametric optimisation into a parametric one, and if set inappropriately it can lead to the resulting policy being suboptimal. Below we examine the different design choices.

1) *Kernel function:* The kernel function defines prior correlations in different parts of the space and if correctly defined can speed up the process of learning. It ideally should be defined to suit the particular domain. However, there are a number of standard kernel functions that can be used.

As explained in Section IV-A2, the summary space  $\mathcal{C}$  consists of both continuous and discrete variables. We therefore define the kernel function on  $\mathcal{C}$  as the sum<sup>3</sup> of kernels on continuous and discrete spaces:

$$k_{\mathcal{C}}(\mathbf{c}, \mathbf{c}') = k_{\text{cont}(\mathcal{C})}(\mathbf{c}, \mathbf{c}') + k_{\text{disc}(\mathcal{C})}(\mathbf{c}, \mathbf{c}') \quad (20)$$

We investigated two standard kernels on the continuous part of the summary space.

The polynomial kernel function is defined as:

$$k(\mathbf{c}, \mathbf{c}'; \sigma_0, p) = (\langle \mathbf{c}, \mathbf{c}' \rangle + \sigma_0^2)^p, \quad (21)$$

where  $\langle \cdot, \cdot \rangle$  is the dot-product, with hyper-parameters  $\sigma_0 \geq 0$  and  $p > 0$ . In the case where  $p = 1$  this is the linear kernel. Higher order polynomial kernels can lead to better results for small input spaces [23]. However, they are less suitable for large input spaces such as the one used here since the prior variance  $k(\mathbf{c}, \mathbf{c})$  grows rapidly with  $|\mathbf{c}| > 1$  [22].

The Gaussian kernel function<sup>4</sup> is defined as:

$$k(\mathbf{c}, \mathbf{c}'; p, \sigma_k) = p^2 \exp\left(-\frac{\|\mathbf{c} - \mathbf{c}'\|^2}{2\sigma_k^2}\right), \quad (22)$$

where  $\sigma_k$  determines how close the points have to be for the values of the function to be correlated, and  $p$  defines the prior variance at each data point since  $k(\mathbf{c}, \mathbf{c}) = p^2$ . The main advantage of the Gaussian kernel over the polynomial kernel is that it is a dot product of an infinite vector of feature functions [22], which gives it the potential to model covariances better. The Gaussian kernel is a stationary kernel meaning that that the covariances only depend on the distance between two summary states. This is in contrast to the polynomial kernel which is non-stationary since the covariance depends on the part of the space where the two summary states are.

The GP-Sarsa algorithm requires that the kernel function is defined on the summary actions space too. Since this is a small discrete space, the  $\delta$ -kernel was used:

$$k(a, a') = \delta_a(a'). \quad (23)$$

The same function is also used on the discrete part of the summary space.

Fig. 3 compares a linear kernel ( $\sigma_0 = 1$  and  $p = 1$ ) with a Gaussian kernel ( $\sigma_k = 5$  and  $p = 4$ ). The results suggest that the Gaussian kernel is more effective as it results in less variable training. This kernel is therefore used in all subsequent experiments unless stated otherwise.

2) *Residual noise:* The second element of policy design is the choice of parameter  $\sigma$ , the noise of the residual  $\Delta Q$  (see Section III-A). As already noted, this controls how much change in the estimate of the Q-function is expected to take place during the process of learning. It is directly related to the randomness of the return (1) and therefore depends on the reward function. The rule of thumb is that  $\sigma$  should be the square root of the half length of the interval over possible values of the reward function. The intuition behind this is that the return in principle can take any value between the highest

<sup>3</sup>A product is also possible but has not been examined here.

<sup>4</sup>Also called the squared exponential kernel function in the literature.

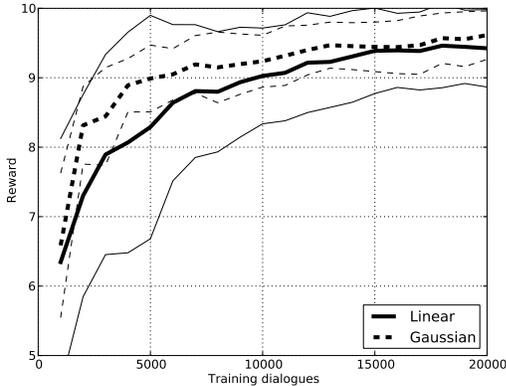


Fig. 3. Influence of the kernel choice on the performance. The performance is measured as the mean average reward as a function of the number of training dialogues. The thinner lines denote one standard error.

and the lowest possible reward. Since the process of learning is non-stationary, the return can vary significantly during the optimisation process. The reward function used here gives 20 for a successful dialogue less the number of turns, where the user can take up to 30 turns. Therefore, the interval is  $[-30, 20]$  so  $\sigma$  is fixed at 5. To experimentally support this claim, values of  $\sigma = 1$ ,  $\sigma = 5$  and  $\sigma = 10$  were tested and the results are shown in Fig 4 where it can be seen that  $\sigma = 5$  does indeed give the best performance.

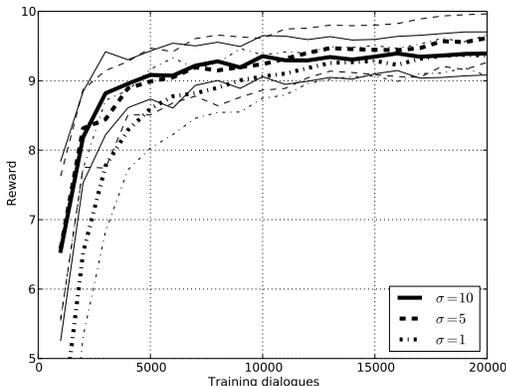


Fig. 4. Influence of the choice of the noise residual parameter  $\sigma$  on the performance. The thinner lines denote one standard error.

3) *Sparsification threshold*: The sparsification threshold  $\nu$  represents a trade-off between the computational complexity and the precision to which the kernel function is computed. Therefore, this should be some small percentage of the minimal value for  $k(\mathbf{c}, \mathbf{c}')$ . Fig. 5 compares performance for  $\nu = 0.1$ ,  $\nu = 1$  and  $\nu = 10$ . As can be seen, the smaller  $\nu$  is the better performance is. It is interesting to note however that the difference in performance between  $\nu = 0.1$  and  $\nu = 1$  is almost negligible in comparison to the difference between  $\nu = 1$  and  $\nu = 10$ . Fig. 6 gives the number of dictionary points for different thresholds. Comparing the two, it can be seen that for a very small improvement in performance the number of dictionary points dramatically increases.

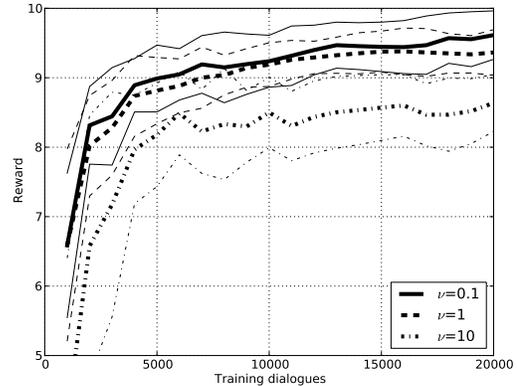


Fig. 5. Influence of the sparsification threshold  $\nu$  on the performance. The thinner lines denote one standard error.

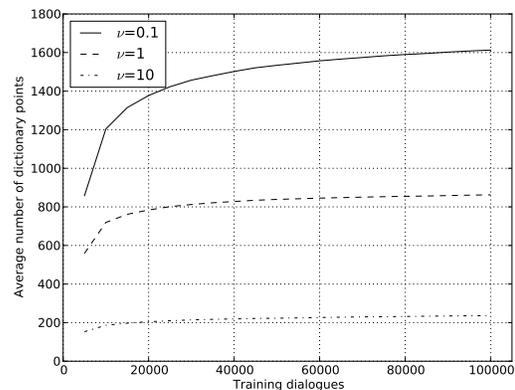


Fig. 6. Average number of dictionary points for different settings of  $\nu$  as a function of the number of training dialogues.

4) *Policy model*: As already noted in Section III-B, the policy model can be defined  $\epsilon$ -greedily (11), using active learning (12), or stochastically (13). It has already been shown that active learning has the potential to lead to faster learning [24] as well as other learning strategies that exploit the variance [36], [37]. Here we focus on the stochastic policy since it avoids the need to explicitly control the exploration rate during training, and it provides a uniform mechanism for both optimisation and normal run-time operation.

Gaussian process estimation has an inherent problem in that the estimate of the variance depends only on the number of points from the input space rather than their values in the output space. This is exacerbated in the set-up used here since the form of Gaussian process approximation used can lead to overly confident estimates [27] whilst at the same time, as input spaces become larger, the same points need to be visited more and more times in order to achieve optimal performance. Hence, when defining a stochastic policy for a large input space, such as the one used here, it is helpful to scale the variances during training to ensure an adequate degree of exploration:  $Q(\mathbf{b}, a) \sim \mathcal{N}(\bar{Q}(\mathbf{b}, a), \eta^2 \text{cov}((\mathbf{b}, a), (\mathbf{b}, a)))$  where  $\eta$  is the scaling factor.

Figs. 7 and 8 show performance during training and testing, respectively, for a stochastic policy with variances scaled by

$\eta = 1$ ,  $\eta = 2$  and  $\eta = 3$ .<sup>5</sup> For comparison, the performance of an  $\epsilon$ -greedy policy is also shown as a baseline.

As can be seen, varying the scale factor  $\eta$  has a significant influence on performance. Setting  $\eta$  to a low value results in the system learning quickly, but also converging to a local optimum. Whilst the best stochastic policy model shows only a small improvement in comparison to  $\epsilon$ -greedy learning during testing (Fig. 8), the performance during training is significantly higher (Fig. 7) because  $\epsilon$ -greedy learning requires an action to be chosen randomly some percentage of time, whereas the stochastic policy always follows the current estimate of the  $Q$ -function. This is particularly important when training in interaction with real users since poor performance during the early stages of training may impact on the user’s willingness to continue to interact with the system.

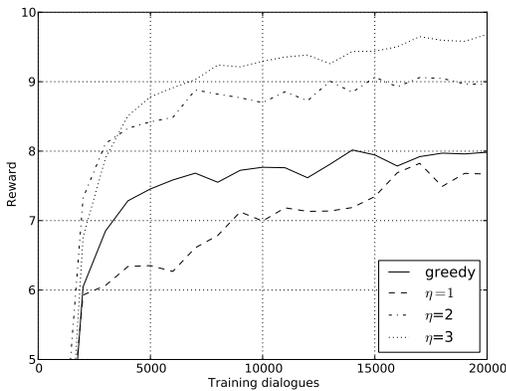


Fig. 7. Influence of variance scaling and the choice of the policy model on the performance during training.

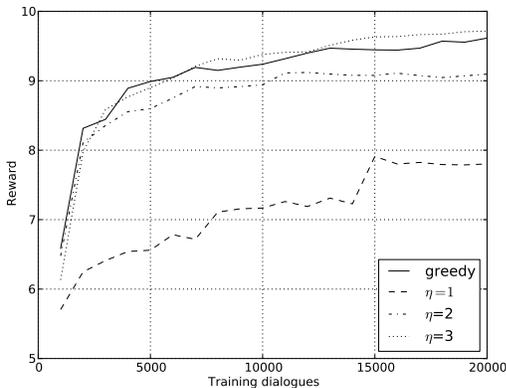


Fig. 8. Influence of variance scaling and the choice of the policy model on the performance during testing.

C. Fast policy optimisation

One of the claimed advantages of the GP-Sarsa algorithm is its ability to learn from a small amount of data. It has been shown that GP-Sarsa outperforms standard non-parametric

<sup>5</sup>Higher values of  $\eta$  did not produce any further improvements in performance.

policy approaches [24]. In this section, we compare the GP-Sarsa learning rate with the NAC gradient-based algorithm.

1) Comparison with standard methods: NAC policy optimisation for the BUDS dialogue manager relies on the parametric representation of the policy  $\pi$  using a set of feature functions  $\mathcal{F}$ , as explained in Section IV-A1. The GP-Sarsa algorithm directly optimises the policy in the summary state  $\mathcal{C}$ .

To measure the rate at which each algorithm learns, the partially optimised policies were evaluated on 1000 random dialogues after being trained on 5000 dialogues. This process was repeated until the policies had been trained with 100,000 dialogues in total. The confusion rate was set to 15% throughout. The simulated user had a different random seed in testing mode to the one set in training mode. The average reward, the success rate and the average number of turns are shown in Figs. 9, 10 and 11 respectively. The error bars represent one standard error.

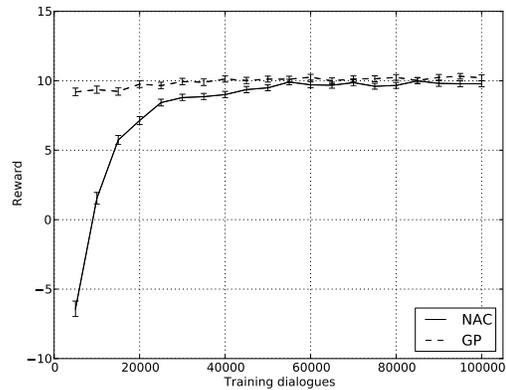


Fig. 9. Comparison of the performance of GP-Sarsa vs the natural actor critic algorithm, where the performance is measured as the average reward as a function of the number of training dialogues.

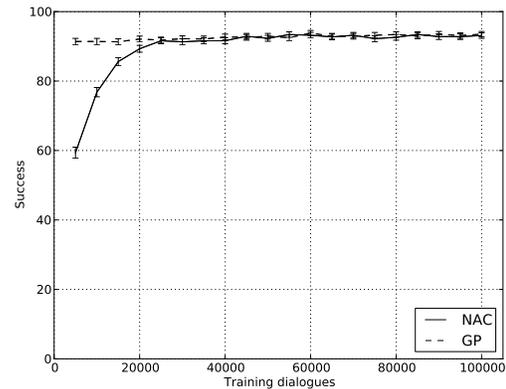


Fig. 10. Comparison of the performance of GP-Sarsa vs the natural actor critic algorithm, where the performance is measured as the average success as a function of the number of training dialogues.

The results show that the GP-based policy reaches a performance comparable to the fully trained NAC policy in only 10,000 dialogues. This is consistent with related research which has shown that Gaussian processes can be deployed to guide gradient descent and improve its efficiency [38].

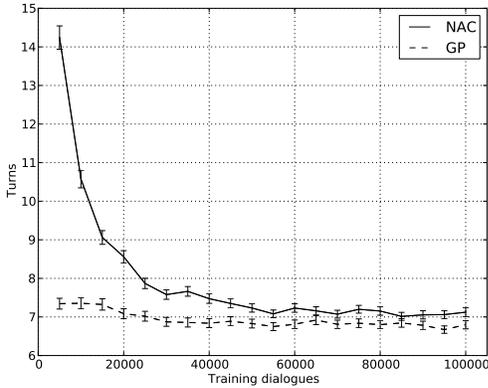


Fig. 11. Comparison of the performance of GP-Sarsa vs the natural actor critic algorithm, where the performance is measured as the average number of turns as a function of the number of training dialogues.

Also, the fully trained GP policy outperforms the NAC policy in terms of dialogue length showing that parametric policy modelling limits the optimality of the solution.

The fully trained GP and NAC policies were also compared by performing 1000 simulated dialogues over a range of confusion rates. The average reward is given in Fig. 12.

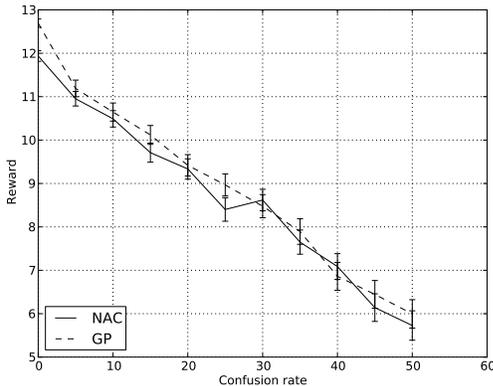


Fig. 12. Comparison of the performance of GP-Sarsa vs the natural actor critic algorithm, where the performance is measured as the average reward as a function of the confusion rate in the user input.

The results show that the GP policy outperforms the NAC policy across confusion rates in the neighbourhood of 15% confusion rate where the policies are trained. This shows that while the GP has a potential for more robust performance, it is important that the training conditions match testing conditions.

Evaluating policy performance on the same simulated user as used for training may be misleading. Hence, the policies were also tested with human users recruited using the Amazon MTurk service. Four policies were compared: partially trained GP and NAC policies, using only 10,000 dialogues, and the fully trained GP and NAC policies.

The results are given in Table I, where for each policy the number of dialogues is given together with the average reward, the success rate, the average number of turns as well as their respective one-standard errors, and the number of dictionary points used by the GP sparsification algorithm.

The average word error rate was 19%. These results show that the partially trained GP policy outperforms the partially trained NAC policy. Indeed, the partially trained GP policy outperforms all of the other policies including the fully trained GP policy. The latter result may come from exploitation of traits in the simulated user that are not normally present in real human interaction.

2) *Real world application:* GP-Sarsa has been previously used for training a dialogue manager in direct interaction with real users [37], however, the summary state space used in that experiment was only four dimensional. In addition, the learning process exhibited inconsistencies due to the users not providing accurate ratings. Here we report on an on-line learning experiment with the Bayesian Update of Dialogue State manager using a refined reward function.

The GP-Sarsa algorithm was implemented in a live telephone-based spoken dialogue system based on the BUDS framework operating with the summary space described in Section IV-A2. The GP-Sarsa configuration consisted of the Gaussian kernel ( $\sigma_k = 5$  and  $p = 4$ ) with a stochastic policy model where  $\eta = 3$ , the noise of the residual  $\sigma = 5$  and the sparsification threshold  $\nu = 0.1$ .

Human users were assigned specific tasks in the Cambridge restaurant domain using the Amazon MTurk service. At the end of each call, users were asked to press 1 if they were satisfied and 0 otherwise. Previous experience suggests that this subjective user feedback is not sufficiently robust to rely on solely for the reward function. Therefore, at the end of each call the objective success was also calculated by comparing the predefined task with the system actions. The user rating was then used to compute the reward function only if it agreed with the objective measure of success. It has been shown that this can result in better policy performance [37].

The performance achieved during on-line learning on the initial 1274 dialogues was compared to the performance of the policy trained on 10,000 dialogues using the simulated user (GPPartlyTrained policy from the previous section). The results are given in Fig. 13 where the moving average success is given for a window of 400 dialogues along with the 95% confidence interval. Of the 1274 dialogues, 1081 had the same subjective and objective success ratings and were used for training. It can be seen from Fig. 13 that the policy trained on-line reaches a comparable asymptotic reward as the simulator-trained policy but in substantially fewer training dialogues. Moreover, the learning curve is smooth and it does not exhibit inconsistencies as was the case in earlier experiments [37] due to the refined definition of the reward function.

Finally, the policy trained on 1000 dialogues with real users was tested under the same conditions as for the simulator trained policies (i.e. the scaling factor  $\eta = 1$ ). Comparing Table II with Table I, it can be seen that the policy trained on real users significantly outperforms the simulator trained policies.

#### D. Policy optimisation on the full belief space

The second major benefit of using Gaussian process policy optimisation is its potential to substantially reduce designer effort in defining the policy representation by avoiding the

TABLE I  
HUMAN EVALUATION OF SIMULATOR TRAINED POLICIES

	#of dialogues	Reward	Success	Turns	#of dict points
NAC-PartlyTrained	402	8.9 ± 0.4	87.3 ± 1.7	8.5 ± 0.2	–
NAC-FullyTrained	415	11.9 ± 0.3	91.8 ± 1.3	6.5 ± 0.2	–
GP-PartlyTrained	400	12.5 ± 0.3	93.5 ± 1.2	6.2 ± 0.2	1569
GP-FullyTrained	397	11.6 ± 0.4	91.2 ± 1.4	6.6 ± 0.2	1699

TABLE II  
HUMAN EVALUATION OF THE POLICY TRAINED ONLINE

	#of dialogues	Reward	Success	Turns	#of dict points
GP-OnlineTrained	410	13.4 ± 0.3	96.8 ± 0.9	6.0 ± 0.1	1132

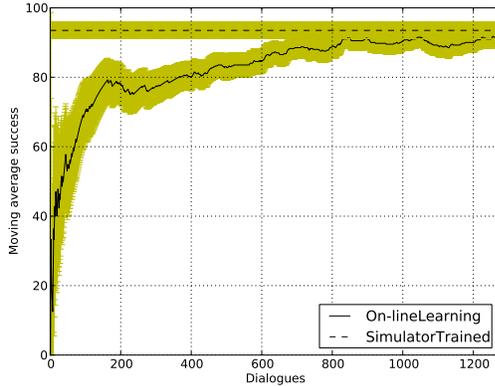


Fig. 13. Comparison of the performance of the policy trained with the simulated user and tested with real people with the policy that is optimised online with interaction with real people. The performance is measured as the moving average success as a function of the number of dialogues.

need to construct *ad hoc* feature mapping functions. In this section, the ability to optimise a policy on the full belief space is examined in practice.

1) *The kernel function on the full belief space:* In the case of the full belief space, the kernel function must be defined over distributions. For the experiments described here, the kernel functions were defined following the approaches described in [39], [40] but modified to suit the dialogue system application. The kernel function was constructed from the sum of individual kernels over the distributions  $\mathbf{b}_k$  for each hidden node  $k$  in the Bayesian network belief state  $\mathbf{b}$ .

We examined three kernel functions. The first is the expected likelihood kernel, which is also a simple linear inner product:

$$k_{\mathcal{B}}(\mathbf{b}, \mathbf{b}') = \sum_k \langle \mathbf{b}_k, \mathbf{b}'_k \rangle, \quad (24)$$

the expected value of distribution  $\mathbf{b}_k$  under distribution  $\mathbf{b}'_k$ .

The second is the Bhattacharyya kernel

$$k_{\mathcal{B}}(\mathbf{b}, \mathbf{b}') = \sum_k \sum_{i=0}^{i < d_k} \sqrt{\mathbf{b}_k(i)} \sqrt{\mathbf{b}'_k(i)}, \quad (25)$$

where  $d_k$  is the dimensionality of the  $k$  concept in the Bayesian network. It is related to Hellinger's distance, which can be viewed as a symmetric approximation of the Kullback-

Leibler divergence [39].

The third kernel that we compare is given by

$$k_{\mathcal{B}}(\mathbf{b}, \mathbf{b}') = \frac{-1}{\log(2)} \sum_k \sum_{i=0}^{i < d_k} \mathbf{b}_k(i) \log\left(\frac{\mathbf{b}_k(i)}{\mathbf{b}_k(i) + \mathbf{b}'_k(i)}\right) + \mathbf{b}'_k(i) \log\left(\frac{\mathbf{b}'_k(i)}{\mathbf{b}_k(i) + \mathbf{b}'_k(i)}\right). \quad (26)$$

This is related to Jansen-Shannon divergence, a symmetric and smoothed variant of the Kullback-Leibler divergence [40].

The kernel functions of the history nodes are defined directly over their respective distributions. Whilst it is possible to calculate the kernel function for the goal nodes in the same way, in this case, the choice of system action, such as *confirm* or *inform*, does not depend on the distribution of actual values, rather it depends on the shape of the distribution. Therefore, the kernel functions for goal nodes are calculated over vectors, where each vector represents the corresponding distribution sorted into order of probability. The only exceptions are the *method* goal node and the *discourse act* node. The former defines whether the user is searching for a venue *by name* or *by constraints* and the latter defines which discourse act the user used, eg. *acknowledgement*, *thank you*. Their kernels are calculated in the same way as for the history nodes.

The sparsification threshold  $\nu$  for this space is set at 0.001. Some intermediate experimentation showed that the kernel requires a finer span, so the threshold here is low. Since the same reward function is used for both the summary space and the full belief space, the variance of the  $Q$ -function residual  $\sigma$  was also fixed at 5. The sampling scale  $\eta$  was set to 3.

2) *Training and evaluation on simulated user:* The GP-Sarsa stochastic policy (13) that operates on the summary space was compared to the two policies operating on the full belief space that use different kernel functions in the same training schedule as described in Section IV-A4.

The results are shown in Fig. 14. As can be seen, the GP policies trained on the full belief space do converge to an optimum, with the Bhattacharyya kernel performing best. Furthermore, increasing the dimensionality of the input space from the summary space to the full belief space did not slow down the process of learning. Overall performance is not however competitive with the summary space stochastic policy. This is disappointing but perhaps not surprising since defining a feature based summary space mapping can be viewed as designing a special kind of kernel that is hand-

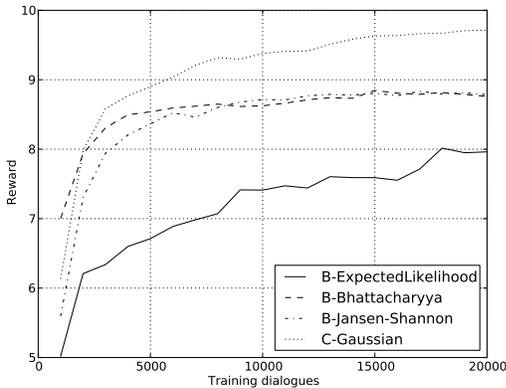


Fig. 14. Comparison of the performance of 3 different policies operating on the full belief space ( $\mathcal{B}$ ) with a policy operating on the summary space ( $\mathcal{C}$ ).

crafted to suit this task well. This is also the first attempt at designing a general-purpose kernel to fit an arbitrarily large belief space, and further work is needed to refine this process.

## V. CONCLUSION

This paper has described how Gaussian processes can be successfully applied to POMDP-based dialogue management. Various aspects of Gaussian process reinforcement learning have been discussed, including sparse approximations, the use of a GP-based stochastic policy model and the GP-Sarsa algorithm.

Using a system based on the Bayesian Update of Dialogue State (BUDS) dialogue manager, various issues have been investigated experimentally both in simulation and with real users. The principal results are firstly that GP-based policy optimisation is faster than conventional gradient based methods and furthermore, performance in the early stages of learning becomes acceptable very quickly making it feasible to train systems from scratch directly in interaction with real users. Secondly, policies trained with real users significantly outperform policies trained on user simulators. Thirdly, although there is clearly more work to be done, it is feasible to construct a GP-based policy directly on the full POMDP belief space obviating the need for hand-crafting a summary space mapping. These three results represent an important step towards building flexible dialogue systems that can range over large and varied domains dynamically adapting to the local context.

Future research in this area will need to address a number of issues. Firstly, methods are needed for optimising the kernel function parameters. It has been shown on a toy dialogue problem that these parameters can be learnt off-line directly from data by maximising the marginal likelihood (9) [24], [23]. This type of approach needs to be extended to work effectively in a full-scale real-world dialogue manager.

Secondly, the work presented here still requires the use of a summary action space and heuristic mapping functions to generate actual system responses. For the experiments here, a simple  $\delta$ -kernel was used as the kernel function over the summary action space. However, Gaussian process reinforcement

learning should allow more elaborate kernel functions to be defined which can be applied directly to the full action space. Since the full action space is not finite, sampling methods would then need to be developed to select an appropriate system response as in [19].

Finally, more work is needed to develop effective training strategies which allow optimisation of GP-Sarsa policies operating on the full belief state space in direct interaction with real users. The preliminary experiments reported in this paper suggest this is possible, but work is clearly needed to achieve performance levels which are competitive with current summary space models. A related issue concerns the design of reward functions, and robust methods of computing them on-line. In all of the work described in this paper, a simple fixed reward for success is given at the end of the dialogue in conjunction with a per turn penalty and the primary measure of success is determined by asking the user to give feedback. This learning regime would not be possible in real applications. Ideally the reward needs to be computed without the conscious engagement of users and it needs to be refined to encourage the type of desirable dialogue behaviours that are embodied in the design of current hand-crafted spoken dialogue systems.

## ACKNOWLEDGMENTS

The authors would like to thank Blaise Thomson, who built the Bayesian Update of Dialogue State system used here, Jost Schatzmann and Simon Keizer, who built the user simulator and Pirros Tsiakoulis for his help with the user trials. The authors would also like to thank Matt Henderson, Kai Yu, Martin Szummer, Catherine Breslin and Dongho Kim for useful comments and discussions.

## REFERENCES

- [1] S. Young, "Talking to Machines (Statistically Speaking)," in *Proc. of ICSLP*, 2002.
- [2] N. Roy, J. Pineau, and S. Thrun, "Spoken dialogue management using probabilistic reasoning," in *Proc. of ACL*, 2000.
- [3] B. Zhang, Q. Cai, J. Mao, E. Chang, and B. Guo, "Spoken Dialogue Management as Planning and Acting under Uncertainty," in *Proc. of Eurospeech*, 2001.
- [4] J. Williams and S. Young, "Partially Observable Markov Decision Processes for Spoken Dialog Systems," *Computer Speech and Language*, vol. 21, no. 2, pp. 393–422, 2007.
- [5] B. Thomson, "Statistical methods for spoken dialogue management," Ph.D. dissertation, University of Cambridge, 2009.
- [6] S. Young, M. Gašić, S. Keizer, F. Mairesse, J. Schatzmann, B. Thomson, and K. Yu, "The Hidden Information State model: A practical framework for POMDP-based spoken dialogue management," *Computer Speech and Language*, vol. 24, no. 2, pp. 150–174, 2010.
- [7] M. Gašić, S. Keizer, F. Mairesse, J. Schatzmann, B. Thomson, K. Yu, and S. Young, "Training and evaluation of the HIS-POMDP dialogue system in noise," in *Proc. of SIGDIAL*, 2008.
- [8] F. Jurčiček, B. Thomson, and S. Young, "Natural actor and belief critic: Reinforcement algorithm for learning parameters of dialogue systems modelled as POMDPs," *ACM Transactions on Speech and Language Processing*, 2011.
- [9] L. Kaelbling, M. Littman, and A. Cassandra, "Planning and Acting in Partially Observable Stochastic Domains," *Artificial Intelligence*, vol. 101, pp. 99–134, 1998.
- [10] J. Pineau, G. Gordon, and S. Thrun, "Point-based value iteration: An anytime algorithm for POMDPs," in *Proc. of IJCAI*, 2003, pp. 1025–1032.
- [11] J. Williams and S. Young, "Scaling POMDPs for Spoken Dialog Management," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 7, pp. 2116–2129, 2007.

- [12] J. Henderson, O. Lemon, and K. Georgila, "Hybrid Reinforcement/Supervised Learning for Dialogue Policies from fixed data sets," *Computational Linguistics*, vol. 34, no. 4, pp. 487–511, 2008.
- [13] B. Thomson and S. Young, "Bayesian update of dialogue state: A POMDP framework for spoken dialogue systems," *Computer Speech and Language*, vol. 24, no. 4, pp. 562–588, 2010.
- [14] L. Li, J. Williams, and S. Balakrishnan, "Reinforcement Learning for Dialog Management using Least-Squares Policy Iteration and Fast Feature Selection," in *Proc. of Interspeech*, 2009.
- [15] P. Crook and O. Lemon, "Lossless Value Directed Compression of Complex User Goal States for Statistical Spoken Dialogue Systems," in *Proc. of Interspeech*, 2011.
- [16] L. Daubigny, M. Geist, and O. Pietquin, "Off-policy Learning in Large-scale POMDP-based Dialogue Systems," in *Proc. of ICASSP*, 2012.
- [17] Amazon, "Amazon Mechanical Turk," 2011, <https://www.mturk.com/mturk/welcome>.
- [18] Y. Engel, S. Mannor, and R. Meir, "Bayes Meets Bellman: The Gaussian Process Approach to Temporal Difference Learning," in *Proc. of ICML*, 2003.
- [19] —, "Reinforcement learning with Gaussian processes," in *Proc. of ICML*, 2005.
- [20] C. E. Rasmussen and M. Kuss, "Gaussian processes in reinforcement learning," in *Advances in Neural Information Processing Systems*, vol. 16. MIT Press, 2004, pp. 751–759.
- [21] M. Deisenroth, C. Rasmussen, and J. Peters, "Gaussian Process Dynamic Programming," *Neurocomputing*, vol. 72, no. 7-9, pp. 1508–1524, 2009.
- [22] C. Rasmussen and C. Williams, *Gaussian Processes for Machine Learning*. Cambridge, Massachusetts: MIT Press, 2005.
- [23] M. Gašić, "Statistical Dialogue Modelling," PhD thesis, University of Cambridge, 2011.
- [24] M. Gašić, F. Jurčićek, S. Keizer, F. Mairesse, J. Schatzmann, B. Thomson, K. Yu, and S. Young, "Gaussian processes for fast policy optimisation of pomdp-based dialogue managers," in *Proc. of SIGDIAL*, 2010.
- [25] D. Cohn, L. Atlas, and R. Ladner, "Improving Generalization with Active Learning," *Machine Learning*, vol. 15, pp. 201–221, 1994.
- [26] MacKay, D.J.C., "Information-based objective functions for active data selection," *Neural Computation*, vol. 4, no. 4, pp. 590–604, 1992.
- [27] J. Quinero-Candela and C. Rasmussen, "A Unifying View of Sparse Approximate Gaussian Process Regression," *Journal of Machine Learning Research*, vol. 6, pp. 1939–1959, 2005.
- [28] Y. Engel, "Algorithms and Representations for Reinforcement Learning," PhD thesis, Hebrew University, 2005.
- [29] D. Nguyen-Tuong and J. Peters, "Incremental sparsification for real-time online model learning," in *Proc. of AISTATS 2010*, 2010.
- [30] B. Scholkopf and A. Smola, *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. Cambridge, MA: MIT Press, 2001.
- [31] J. Peters and S. Schaal, "Natural Actor-Critic," *Neurocomputing*, vol. 71, pp. 1180–1190, 2008.
- [32] J. Schatzmann, "Statistical User and Error Modelling for Spoken Dialogue Systems," Ph.D. dissertation, University of Cambridge, 2008.
- [33] S. Keizer, M. Gašić, F. Jurčićek, F. Mairesse, B. Thomson, K. Yu, and S. Young, "Parameter estimation for agenda-based user simulation," in *Proc. of SIGDIAL*, 2010.
- [34] B. Thomson, M. Gasic, M. Henderson, P. Tsiakoulis, and S. Young, "N-Best error simulation for training spoken dialogue systems," in *Proc. of SLT*, 2012.
- [35] F. Jurčićek, S. Keizer, M. Gašić, F. Mairesse, B. Thomson, K. Yu, and S. Young, "Real user evaluation of spoken dialogue systems using Amazon Mechanical Turk," in *Proc. of Interspeech*, 2011.
- [36] L. Daubigny, M. Gašić, S. Chandramohan, M. Geist, O. Pietquin, and S. Young, "Uncertainty management for on-line optimisation of a POMDP-based large-scale spoken dialogue system," in *Proc. of Interspeech*, 2011.
- [37] M. Gašić, F. Jurčićek, B. Thomson, K. Yu, and S. Young, "On-line policy optimisation of spoken dialogue systems via live interaction with human subjects," in *Proc. of ASRU*, 2011.
- [38] H. Jakab and L. Csató, "Reinforcement learning with guided policy search using gaussian processes," in *Neural Networks (IJCNN), The 2012 International Joint Conference on*. IEEE, 2012, pp. 1–8.
- [39] T. Jebara, R. Kondor, and A. Howard, "Probability product kernels," *J. Mach. Learn. Res.*, vol. 5, pp. 819–844, Dec. 2004.
- [40] M. Hein and O. Bousquet, "Hilbertian metrics and positive definite kernels on probability measures," in *Proc. of AISTATS 2005*, 2004.



organising committee for Young Researcher's Roundtable on Spoken Dialogue Systems in 2009.

**Milica Gašić** Milica Gašić is a Research Associate in the Dialogue Systems Group at the University of Cambridge. She graduated in Computer Science and Mathematics from the University of Belgrade in 2006. Since then she has been at the University of Cambridge first as an MPhil student reading Computer Speech, Text and Internet Technology. Subsequently, she did a PhD in Statistical Dialogue Modelling in 2011. She has published around 20 peer-reviewed conference and journal papers in the area of dialogue management. She served on the



Senior Pro-Vice-Chancellor.

**Steve Young** Steve Young received a BA in Electrical Sciences from Cambridge University in 1973 and a PhD in Speech Processing in 1978. He was elected to the Chair of Information Engineering at Cambridge University in 1994. He was a co-founder and Technical Director of Entropic Ltd from 1995 until 1999 when the company was taken over by Microsoft. After a short period as an Architect at Microsoft, he returned full-time to Cambridge University in January 2001 where he was Head of Information Engineering until 2009 and is now

His research interests include speech recognition, expressive synthesis and spoken dialogue systems. He has written and edited books on software engineering and speech processing, and he has published as author and co-author, more than 200 papers in these areas. He is a Fellow of the UK Royal Academy of Engineering, the IET, IEEE and RSA. He served as the senior editor of *Computer Speech and Language* from 1993 to 2004, and he was Chair of the IEEE Speech and Language Technical Committee from 2008 to 2010. He was the recipient of an IEEE Signal Processing Society Technical Achievement Award in 2004. He was elected a Fellow of the International Speech Communication Association (ISCA) in 2009 and he was the recipient of the ISCA Medal in 2010 for Scientific Achievement.