



HAL
open science

Bayesian update of dialogue state: A POMDP framework for spoken dialogue systems

Blaise Thomson, Steve Young

► **To cite this version:**

Blaise Thomson, Steve Young. Bayesian update of dialogue state: A POMDP framework for spoken dialogue systems. *Computer Speech and Language*, Elsevier, 2010, 24 (4), pp.562. 10.1016/j.csl.2009.07.003 . hal-00621617

HAL Id: hal-00621617

<https://hal.archives-ouvertes.fr/hal-00621617>

Submitted on 11 Sep 2011

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Accepted Manuscript

Bayesian update of dialogue state: A POMDP framework for spoken dialogue systems

Blaise Thomson, Steve Young

PII: S0885-2308(09)00049-7
DOI: [10.1016/j.csl.2009.07.003](https://doi.org/10.1016/j.csl.2009.07.003)
Reference: YCSLA 428

To appear in: *Computer Speech and Language*

Received Date: 17 October 2008
Revised Date: 19 June 2009
Accepted Date: 25 July 2009



Please cite this article as: Thomson, B., Young, S., Bayesian update of dialogue state: A POMDP framework for spoken dialogue systems, *Computer Speech and Language* (2009), doi: [10.1016/j.csl.2009.07.003](https://doi.org/10.1016/j.csl.2009.07.003)

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Bayesian update of dialogue state: A POMDP framework for spoken dialogue systems

Blaise Thomson and Steve Young

University of Cambridge

Abstract

Based on the framework of partially observable Markov decision processes (POMDPs), this paper describes a practical real-time spoken dialogue system in which the underlying belief state is represented by a dynamic Bayesian Network and the policy is parameterized using a set of action-dependent basis functions. Tractable real-time Bayesian belief updating is made possible using a novel form of Loopy Belief Propagation and policy optimisation is performed using an episodic Natural Actor Critic algorithm. Details of these algorithms are provided along with evaluations of their accuracy and efficiency.

The proposed POMDP-based architecture was tested using both simulations and a user trial. Both indicated that the incorporation of Bayesian belief updating significantly increases robustness to noise compared to traditional dialogue state estimation approaches. Furthermore, policy learning worked effectively and the learned policy outperformed all others on simulations. In user trials the learned policy was also competitive, although its optimality was less conclusive. Overall, the Bayesian update of dialogue state framework was shown to be a feasible and effective approach to building real-world POMDP-based dialogue systems.

Key words: Dialogue Systems, Robustness, POMDP, Reinforcement Learning

1 Introduction

The ability to converse is often considered a defining characteristic of intelligent behaviour. Human-machine dialogue therefore has great significance both academically and as a practical method of communicating with machines for commercial applications. Evidence for the latter can be found in the frequent

Email address: {brmt2,sjy}@eng.cam.ac.uk (Blaise Thomson and Steve Young).

1 use of Spoken Dialogue Systems (SDSs) in applications such as access to in-
2 formation and services, interaction with robots, playing games and provid-
3 ing customer support. Any improvement in the effectiveness of these systems
4 would have far reaching implications.
5

6 Commercial dialogue systems are typically implemented by flowcharting sys-
7 tem prompts along with possible user responses (Pieraccini and Huerta, 2008).
8 The system is represented as a graph, sometimes called the *call flow*, where
9 nodes represent prompts or actions to be taken by the system and the arcs
10 give the possible responses. Formally, the system is therefore a Finite State
11 Automaton (FSA). This approach has been effective commercially because it
12 is straightforward to implement and system prompts can be designed to elicit
13 highly restricted user responses. However, the resulting dialogues can become
14 frustrating for users as their choice is severely limited. Further problems arise
15 when speech recognition and understanding errors occur. In some cases, the
16 system might accept information that is in fact incorrect and elaborate error
17 detection and correction schemes are then required to recover from this. As
18 a result, commercial SDSs are seldom robust to high noise levels and require
19 significant effort and cost to develop and maintain.
20
21
22
23
24

25 Researchers have attempted to overcome these issues of robustness and devel-
26 opment cost in various ways, many of which fall into two categories. The first is
27 to model dialogue as a conversational game with rewards for successful perfor-
28 mance. Optimal action choices may then be calculated automatically, reducing
29 the costs of maintenance and design as well as increasing performance. These
30 models have largely been based on the Markov Decision Process (MDP), as
31 for example in Levin et al. (2000); Walker (2000); Scheffler (2002); Pietquin
32 (2004) and Lemon et al. (2006). The choice of actions for a given internal
33 system state is known as the *policy* and it is this policy which MDP models
34 attempt to optimise.
35
36
37
38

39 A second avenue of research has been to use a statistical approach to modelling
40 the uncertainty in dialogue. This allows for simpler and more effective methods
41 of dealing with errors. Pulman (1996); Horvitz and Paek (1999) and Meng et al.
42 (2003) all suggest early forms of this approach. These statistical systems view
43 the internal system state as representing a set of *beliefs* about the state of the
44 environment. The true state of the environment is hidden and must be inferred
45 from observations, usually via Bayesian inference. The environment state is
46 often separated into different components, called *concepts*, each of which has
47 a set of possible *concept values*. Sometimes these concepts are called *slots*.
48
49
50
51

52 More recently, there have been several attempts to combine the use of statisti-
53 cal policy learning and statistical models of uncertainty. The resulting frame-
54 work is called the Partially Observable Markov Decision Process (POMDP)
55 (Roy et al., 2000; Williams and Young, 2006b; Bui et al., 2007). Figure 1 shows
56
57
58
59
60
61
62
63
64
65

how this framework compares to those suggested previously. The attention of this paper is focussed on the lower half of this figure.

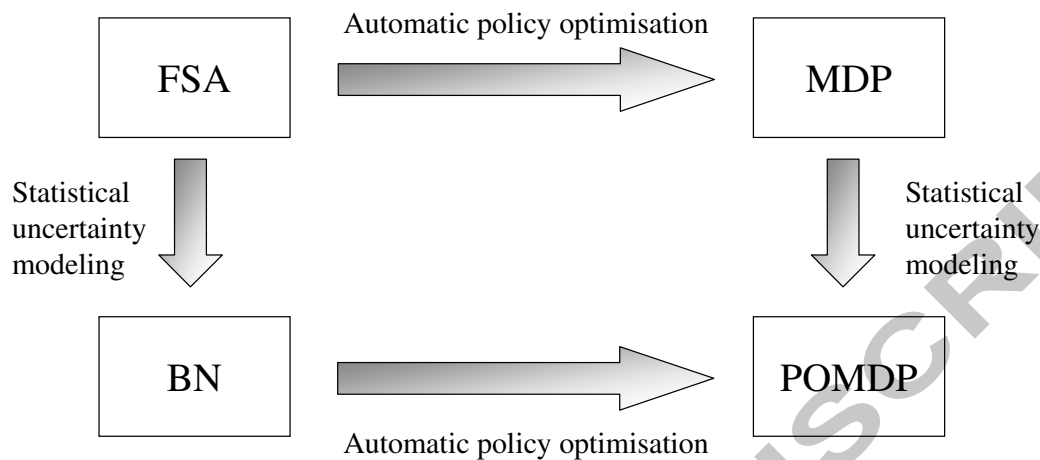


Fig. 1. Frameworks for modelling uncertainty and policy optimisation in Spoken Dialogue Systems: Finite State Automata (FSA), Bayesian Networks (BN), Markov Decision Processes (MDP) and Partially Observable MDPs (POMDP).

Tractability is a major concern whenever a statistical model of uncertainty is used since statistical models require that a probability be associated with every possible state of the environment. Unless significant assumptions and approximations are taken, the result is intractable for any real-world system. Many authors suggest using Bayesian network algorithms as a solution (Pulman, 1996; Horvitz and Paek, 1999; Meng et al., 2003; Bui et al., 2007). Young et al. (2007) take a different approach, grouping the state space into *partitions* where states have similar characteristics and pruning unlikely cases. Another alternative based on particle filters is suggested by Williams (2007b). However, all of these models are restricted to cases where the user's goal stays constant and there are limited dependencies between concepts or there are limited concepts or concept values.

This paper presents a Bayesian Network-based architecture which allows the user goal, user input act and dialogue history to be modelled in a flexible and scalable manner. To enable tractable belief updating, a novel form of the well-known Loopy Belief Propagation (LBP) algorithm (Bishop, 2006, Ch. 8) has been developed which uses partitioning of variable value-sets and tying of transition probabilities to achieve significant efficiency gains.

Previous attempts to build POMDP-based spoken dialogue systems have relied on mapping the state space into a much more compact summary space in order ensure tractable policy representation and optimisation (Williams and Young, 2007; Thomson et al., 2007). A further contribution of this paper is to show how a component-based policy can be defined over the full state space and optimised using the Natural Actor Critic (NAC) algorithm Peters et al.

1 (2005). The NAC algorithm has several advantages over previously proposed
2 approaches. Most importantly, the use of a component-based policy enables
3 the system to learn to differentiate between both slot-level and higher-level
4 actions simultaneously. In contrast, the CSPBVI algorithm learns actions for
5 individual slots and then uses heuristics to decide which action to take. Simi-
6 larly, the summarised form of Q-learning presented in Thomson et al. (2007)
7 places much of the high-level decision making into the heuristic mapping be-
8 tween the summary and the real-world action. Other advantages of the NAC
9 algorithm include its potential for online policy adaptation and its relatively
10 strong convergence properties.
11
12
13

14
15
16
17
18 The paper is organised as follows. Section 2 discusses the representation of
19 dialogue state and associated belief update algorithms. The section starts
20 by introducing a general model of dialogue, POMDPs, Bayesian Networks
21 and other graphical models. Several network structures which are particularly
22 useful for dialogue systems are then described in Section 2.2. Algorithms for
23 implementing dialogue state updates are then discussed in Section 2.3, followed
24 by methods to improve their efficiency in Sections 2.4 and 2.5. An analysis
25 of the efficiency gains achieved by the proposed algorithms is presented in
26 Section 2.6.
27
28
29

30
31
32
33
34 Section 3 is concerned with policy optimisation, focusing on continuous-state
35 Markov Decision Processes which are introduced in Section 3.1. A technique
36 for using domain knowledge to simplify the choice of actions is presented in
37 Section 3.2, followed in Section 3.3 by a discussion of how component-based
38 parametric policies can be used for learning in a dialogue system. Section 3.4
39 then gives a detailed explanation of parameter optimisation using the Natural
40 Actor Critic algorithm. The section ends with a short comparison with other
41 learning algorithms.
42
43
44
45

46
47
48
49 An evaluation of a real-time spoken dialogue system which uses the Bayesian
50 updating framework is presented in Section 4. The evaluation is based on both
51 simulation and a user trial and the results of these are given in sections 4.3
52 and 4.4, respectively. Both the simulations and the user trial indicate that
53 the Bayesian updating framework outperforms standard alternatives. Section
54 5 concludes the paper and gives directions for future research.
55
56
57
58

2 Maintaining Dialogue State

2.1 A theoretical model for dialogue

Current models of dialogue all start from the assumption that a dialogue consists of a sequence of user and system turns. Dialogue management is concerned with deciding what action to take at each system turn and this decision is always based on some form of internal state. The representation of this internal state cannot be a perfect because the system does not have perfect knowledge of the environment. Instead, the internal state must represent the system's *beliefs* about what the user intends and what has happened in the dialogue. The standard notation for the internal state varies between the POMDP and MDP literatures¹. In this paper the system's internal state will be consistently denoted $b \in \mathcal{B}$, and is always called the *belief state*. This includes all different frameworks (including FSA and MDP) since in all cases the internal state is a representation of the system's beliefs. The belief state is often called the *system state* but that will be avoided here to avoid confusion with the *environment state* defined below.

What the system chooses to do during its turn is called the *machine action* or *system action* and is chosen from a predefined set, $m \in \mathcal{M}$. The set of distributions over actions is denoted $\Pi(\mathcal{M})$. Policies, $\pi : \mathcal{B} \rightarrow \Pi(\mathcal{M})$, are then mappings from belief states to probabilistic action choices. The development of the belief state over time depends on the responses the system obtains from its environment. A set of observations, $o \in \mathcal{O}$ is defined, which describes what the system can observe about the world around it². Given the previous belief state, b , the last machine or system action, m , and the last observation, o , the system transitions into a new belief state, b' . The function which defines this is called the *belief state transition function*.

Clearly the actions, belief states and observations are all indexed by the turn number. When it is important to note this dependency they are denoted m_t , b_t and o_t . While the system is in state b_t it will choose action m_t according to the distribution determined by the policy, $\pi(b_t)$. The system then observes observation o_{t+1} and transitions to a new system belief state b_{t+1} . When the

¹ In the MDP literature the internal system state is denoted s . This symbol is reserved for the *environment state* in the POMDP literature.

² In traditional dialogue systems the actions of the system are the system prompts and the observations are the recogniser inputs. More recently, researchers have expanded both the action and observation sets. For example, in a DSL modem troubleshooting task, the actions might also include testing a user's internet connectivity or checking that their password is correct. Observations can include responses from these tests or any other perceptual input from the environment (Williams, 2007a).

time dependence is insignificant, the t is omitted and a prime symbol is used to denote the next time step (e.g. $o' = o_{t+1}$).

The definition of the belief state and its transitions is the key feature which distinguishes models that use a statistical model of uncertainty from those which do not. The essential feature of traditional MDP and FSA models is that the number of belief states is finite. Measures of uncertainty are maintained via hand-crafted rules which represent transitions through varying levels of confidence. Thus, although the internal state can still represent the system's beliefs about its environment, it must be explicitly encoded by the system designer.

In POMDP models the system-state transitions are defined indirectly via probability rules. The model assumes a set of underlying *environment states*, $s \in \mathcal{S}$. Next, a stationary observation function, $P(o_t | s_t, m_{t-1})$, is assumed in which the observation probability is conditionally dependent on only the current environment state, s_t , and the last machine action, m_{t-1} . Changes in the environment state are governed by a constant transition function $P(s_{t+1} | s_t, m_t)$, which assumes that changes in the environment state are dependent only on the last environment state and action but not the full state history. This is called the Markov assumption. The belief state, b_t , is the probability distribution over all possible environment states. Under the above assumptions, this distribution can be calculated using standard probability theory. It is therefore a fully statistical representation of the system's beliefs about its environment.

The difficulty with using a fully statistical approach is that belief updates quickly become intractable. A naive enumeration of all environment states grows exponentially with the number of concepts. Real-world systems must therefore use computationally efficient methods to update the belief distribution over these states.

2.2 Bayesian network models of dialogue

In order to obtain computationally efficient algorithms, the structure of the domain under consideration must be exploited. Bayesian networks provide a graphical representations of statistical models which give an intuitive representation of the structure in a system and also facilitate the use of computationally effective updating algorithms. A detailed introduction can be found in Bishop (2006, Ch. 8).

A Bayesian network is defined to be a directed acyclic graph where the nodes represent random variables, and the arcs represent conditional dependencies. The joint distribution of all variables in the graph factorises as the product of the conditional probability of each variable given its parents in the graph. Fig-

ure 2 gives an example network, representing the assumptions of the POMDP model. Networks which repeat the same structure at each point in time, as in this example, are known as *dynamic Bayesian networks (DBNs)*. The sections of the network corresponding to a particular time are called *time-slices*.

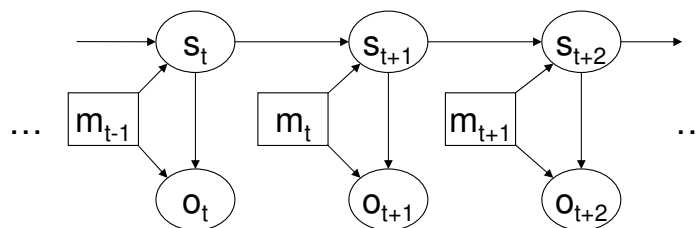


Fig. 2. A portion of a Bayesian Network representing the POMDP model. Decisions are drawn in a rectangle to show that they are actions of the system rather than observed random variables.

The usefulness of Bayesian networks for dialogue comes from allowing further factorisations of the environment state, s_t . When factorising out components of the state, Bayesian network algorithms provide efficient methods for calculating updated marginal distributions of the factors. Updating the beliefs of the system then becomes more tractable.

2.2.1 Structure of dialogue networks

In a dialogue system the environment state, s_t , is highly structured. The exact structure of the Bayesian network used to describe the environment state will depend on the application but there are various structural elements which can often be reused.

One useful factorisation is to separate the environment state into three components: $s_t = (g_t, u_t, h_t)$, where g_t is the long term goal of the user, u_t is the true user act and h_t is the dialogue history (Williams et al., 2005)³. The observed user utterance is then conditionally dependent only on the true user act.

In many systems, further structuring is possible by separating the state into *concepts*, $c \in \mathcal{C}$. In a tourist information system, for example, typical concepts might be the type of “food” served or the “area” in which a venue is located. The state is now factorised into sub-goals, sub-user-acts and sub-histories for each concept. These are denoted $g_t^{(c)}$, $u_t^{(c)}$ and $h_t^{(c)}$. Typically $u_t^{(c)}$ is assumed to depend only on the sub-goal $g_t^{(c)}$. The sub-history $h_t^{(c)}$ will depend on the sub-user-act $u_t^{(c)}$ and the previous sub-history $h_{t-1}^{(c)}$.

It is useful to also maintain a node for the overall user act u_t . To simplify the updates, a one-to-one mapping is defined between the sub-acts,

³ In Williams et al. (2005) the notation is slightly different with $s_u = g$, $a_u = u$ and $s_d = h$.

$u_t^{(c)}$ and this overall act. The sub-acts will represent parts of the action that might be relevant for a particular concept. For example, if the user says “I would like a Chinese restaurant”, the overall act might be represented as “inform(type=restaurant, food=Chinese)”. This will be split so that the sub-act for the *type* concept is “inform(type=restaurant)”, while the sub-act for *food* is “inform(food=Chinese)”. Acts which cannot be split are associated with every concept. An example of the latter is the user saying “Yes”. This will result in an “affirm” act being associated with all slots.

The history node allows the system designer to store sufficient information about the history in order to make coherent decisions. This may be very restricted and in many cases only requires a few values. Indeed, the three states “nothing said”, “mentioned” and “grounded” will normally be adequate.

The dependencies of the sub-goals must be limited to enable tractability. A useful method of implementing this is to add a validity node, v_t , for each concept. This node has a deterministic dependency on its parent (or parents) which decides whether the associated sub-goal is relevant to the overall user goal or not. Validity nodes can only take two values - “Applicable” and “Not Applicable”. If a node is “Not Applicable” then the associated user sub-goal is forced to also be “Not Applicable”. Otherwise the user sub-goal will depend on its previous value with some probability of change. Figure 3 shows the resulting network for a system with two concepts: type (of venue) and food. Note that the user goal and user act are assumed to be conditionally independent of the history.

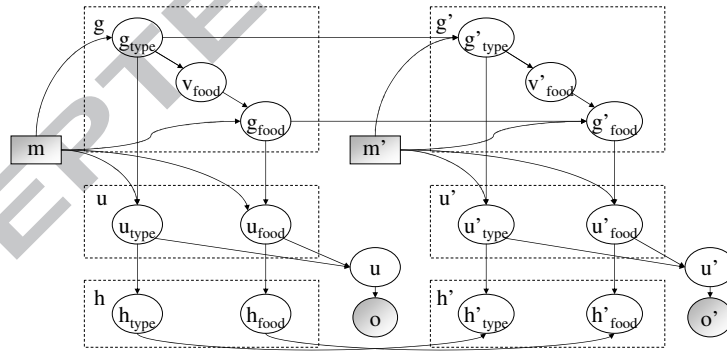


Fig. 3. An example factorisation for the Bayesian network representing part of a tourist information dialogue system. The “type” node represents the type of venue being sought and is assumed to always be relevant so v_{type} is omitted from the diagram.

Validity nodes provide a simple method of switching on relevant parts of the network as the user’s intention is clarified during the dialogue. For example, when a user asks about hotels in a tourist information system, the “food” concept would not be applicable. The concept might, however, be relevant when talking about restaurants in the same system. The food node would

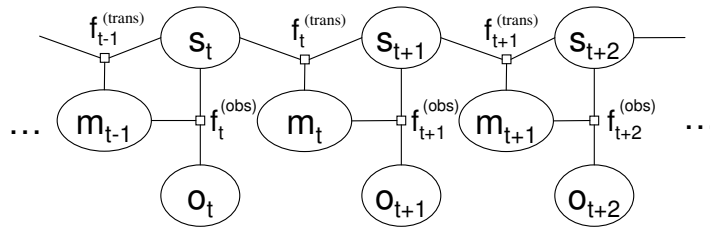
1 therefore be relevant if the user is looking for a restaurant and irrelevant
 2 otherwise. If it becomes clear that the user is looking for a restaurant, the
 3 validity of the food node increases and that part of the network increases in
 4 probability.

5
 6 The remaining nodes in the network will be dependent on the application. It
 7 is important to note that the nodes need not necessarily be *slots* that are filled
 8 with *slot values*. They can represent any random variable in the system⁴.
 9

10 2.3 Loopy belief propagation and factor graphs

11
 12
 13
 14
 15
 16 Once a suitable structure is found to represent the complete environment
 17 state, an efficient algorithm is required to update the beliefs. When performing
 18 calculations, many algorithms are better expressed using a different graphical
 19 representation called a *factor graph* (Kschischang et al., 2001). An example
 20 factor graph is shown in Figure 4 below. Factor graphs are undirected bipartite
 21 graphs, with two types of node. One type represents random variables (drawn
 22 as a circle), while the other type represents factors (drawn as a small square).
 23 The assumption encoded within the graph is that the joint distribution over
 24 all random variables can be written as a product of factor functions, one for
 25 each factor node. These factors are a function of only the random variables
 26 connected to the factor node in the graph.
 27
 28
 29
 30

31 Since a Bayesian network defines a factorisation of the joint distribution
 32 into conditional probability factors, there is a direct mapping from Bayesian
 33 networks to factor graphs. Figure 4 is a factor graph representation of the
 34 POMDP assumptions, previously depicted as a Bayesian network in Figure
 35 2. $f_t^{(trans)}$ represents the environment state's transition function, thus
 36 $f_t^{(trans)}(s_t, s_{t+1}, m_t) = P(s_{t+1}|s_t, m_t)$ and $f_t^{(obs)}$ represents the observation func-
 37 tion, $f_t^{(obs)}(s_t, o_t, m_{t-1}) = P(o_t|s_t, m_{t-1})$.
 38
 39
 40



51 Fig. 4. A portion of a factor graph representing the POMDP model.

52
 53 ⁴ Indeed, as mentioned in footnote 2, Williams (2007a) has described a Bayesian
 54 network structure for a POMDP-based dialogue system to troubleshoot internet
 55 connections in which nodes include environment conditions such as whether there
 56 is “no power” or “no network”.
 57
 58
 59
 60
 61
 62
 63
 64
 65

Algorithm 1 Loopy Belief Propagation

initialize: Set all messages equal to one.

Let $Y = \{\mathbf{x} = (x_1, x_2, \dots, x_N)^\top : x_i \text{ is a possible value of } X_i\}$.

Note that Y enumerates the set of all combinations of the X_i random variables.

repeat

Choose a factor f_a to update. Suppose this is connected to variables X_1, X_2, \dots, X_N .

for each variable, X_i , connected to the factor **do**

(a) Update the messages **out** of the factor.

$$(\forall x'_i) \quad \mu_{f_a \rightarrow X_i}(x'_i) = \sum_{\mathbf{x} \in Y, x_i = x'_i} f_a(\mathbf{x}) \prod_{j \neq i} \mu_{X_j \rightarrow f_a}(x_j). \quad (2)$$

(b) Update the messages **into** nearby factors

for each factor $b \neq a$, connected to X_i **do**

$$(\forall x'_i) \quad \mu_{X_i \rightarrow f_b}(x'_i) = \prod_{c \neq b} \mu_{f_c \rightarrow X_i}(x'_i). \quad (3)$$

end for

end for

until convergence

Factor graphs facilitate the computation of updated marginal distributions given each new set of input observations (i.e. user acts). Since exact computation is generally intractable, approximate methods must be used. Empirically Loopy Belief Propagation (LBP) has been found to be sufficiently accurate and efficient for a variety of applications, and LBP forms the basis of the update algorithms developed here.

Loopy belief propagation maintains a set of messages for each arc in the model. If there is an arc between the random variable X_i and the factor f_a then two messages are defined. $\mu_{X_i \rightarrow f_a}(x_i)$ represents the message from the variable to the factor (into the factor) and $\mu_{f_a \rightarrow X_i}(x_i)$ represents the message from the factor to the variable (out of the factor). Both of these are functions of the possible values of X_i .

The details for calculating these messages are given in Algorithm 1. A derivation of the algorithm is beyond the scope of this paper but may be found in Bishop (2006, Ch. 8). Once the messages have been computed, the marginal probability of any variable, X_i , can be calculated from the messages into that variable from the neighbouring factors, $a \in ne(X_i)$. If k is a normalising constant then:

$$p(x_i) \approx k \prod_{a \in ne(X_i)} \mu_{f_a \rightarrow X_i}(x_i). \quad (1)$$

1 The factor updates in LBP may be completed in any sequence although in
2 practice the system designer will aim to choose a sequence which minimises
3 the number of iterations. The updating should continue until the messages
4 no longer change significantly. At this stage the algorithm is said to have
5 *converged* and the resulting set of messages will constitute a *fixed point* of the
6 algorithm. In cases where the factor graph has a tree structure, the algorithm
7 is exact and will converge after a single breadth-first iteration through the tree
8 followed by a single reverse iteration. On more complex graphs, convergence
9 is not guaranteed, although convergence does usually occur in practice.
10
11
12
13

14 *2.3.1 Limited time-slices*

15
16 One practical difficulty in using the LBP algorithm is that the number of
17 nodes in the network will grow with time. Hence, it may be preferable to
18 limit the number of time-slices, n , that are maintained (Murphy, 2002). The
19 algorithm then proceeds as before except that the factor updates are restricted
20 to the most recent n time-slices. The marginal distributions for variables that
21 are connected to factor nodes in the most recent n time-slices must still be
22 maintained but any other information may be removed.
23
24
25
26

27 In the special case when only one time-slice is maintained and the single
28 slice network has a tree structure, this approach reduces to a special case of
29 the Boyen-Koller algorithm. No iteration will be required because of the tree
30 structure and the error, in terms of expected KL-divergence from the exact
31 updated distributions, will remain bounded over time as proved in Boyen and
32 Koller (1998).
33
34

35 Experiments show that updating the entire history tends to give better accu-
36 racy whereas updating only the most recent time-slice gives faster computa-
37 tion and avoids convergence problems (Murphy, 2002). Further details on the
38 computational efficiency of the algorithms are given in Section 2.6 and the
39 convergence and accuracy issues are discussed further in Appendix B.
40
41
42
43

44 *2.4 Grouped loopy belief propagation*

45
46
47
48 The standard LBP algorithm improves efficiency by exploiting the dependen-
49 cies between concepts. In a dialogue system, a significant issue arises where
50 a concept can take a large number of values. A tourist information system,
51 for example, might have many thousands of options for the name of the avail-
52 able venues. Using loopy belief propagation will help to mitigate the effect of
53 dependencies but updates still become intractable when there are too many
54 concept values.
55
56
57
58
59
60
61
62
63
64
65

One solution to this problem is exploited in the Hidden Information State (HIS) model of Young et al. (2007). Given an observation, many of the possible environment states are indistinguishable and can therefore be grouped. In the HIS system this is done by creating partitions of the true goal space. However, since the partitions must be computed over the full unfactorised goal, it is not possible to exploit the conditional independencies between concepts. The approach discussed here builds on the HIS idea of partitions, but uses it in the context of Bayesian networks and factor graphs. Partitions of concept-values for a particular node are considered rather than partitions of the unfactorised goal. If observations do not distinguish between different concept-values, the messages for those values will be identical and need only be calculated once.

2.4.1 An example

As an example, consider the loopy belief updates of sub-goals such as the food sub-goal in the state factorisation shown in Figure 3. The true value of each sub-goal is $g_t^{(c)}$, while as explained in Section 2.2.1, $v_t^{(c)}$ and m_t represent the validity of the sub-goal and the machine action respectively.

An important assumption about the goal transitions is now made. Given a machine action m_t , there is some probability of staying in the same state, $\theta_{m_t,1}$ and a constant probability, $\theta_{m_t,2}$, of changing states. For the moment, the concept, c , is assumed to be applicable so $v_t^{(c)}$ is ignored. Written as an equation, the factor representing the sub-goal transition function is:

$$f_a(g_{t+1}^{(c)}, g_t^{(c)}, m_t) = P(g_{t+1}^{(c)} | g_t^{(c)}, m_t) = \begin{cases} \theta_{m_t,1} & \text{if } g_t^{(c)} = g_{t+1}^{(c)} \\ \theta_{m_t,2} & \text{if } g_t^{(c)} \neq g_{t+1}^{(c)} \end{cases}. \quad (4)$$

Equation 2 of the LBP algorithm requires a summation over the values of this factor, fixing one value and weighting by the product of the messages from the surrounding variables, of which there is only one in this case. Since the machine action m_t is known, it is also fixed. The resulting updates are as follows:

$$\mu_{f_a \rightarrow g_{t+1}}(g') = \sum_g P(g_{t+1} = g' | g_t = g, m_t) \mu_{g_t \rightarrow f_a}(g) \quad (5)$$

$$\mu_{f_a \rightarrow g_t}(g) = \sum_{g'} P(g_{t+1} = g' | g_t = g, m_t) \mu_{g_{t+1} \rightarrow f_a}(g') \quad (6)$$

where the concept superscript (c) has been omitted for clarity.

These equations require a summation over the entire set of possible user goals.

1 This is clearly inefficient, since when there is nothing to distinguish different
 2 goals, the messages into the factor ($\mu_{g_t \rightarrow f_a}$ and $\mu_{g_{t+1} \rightarrow f_a}$) will be constant.
 3 They can therefore be factored out of the summation, leaving a sum over
 4 the factor values which can be calculated offline. The result will be a set of
 5 equations which only require a summation for each group of user goals for
 6 which the messages are equal.
 7

8 This idea is formalised by partitioning the possible user goals into groups. A
 9 collection of sets, Z_1, Z_2, \dots, Z_N , is defined as a set of partitions of the values
 10 for g_t and g_{t+1} such that whenever $g \in Z_i$ and $g \in Z_j$ then $i = j$ (i.e. there is no
 11 overlap). In the example of the “food” concept, one might have three values
 12 “Chinese”, “Indian” and “French”. When the user has only spoken about
 13 “Chinese” food it is reasonable to split the values into two partitions - one
 14 containing only the “Chinese” option and the other containing the “Indian”
 15 and “French” options. Note that the same partitioning may be used for both
 16 the previous and the next time step since the set of values is the same.
 17
 18
 19
 20

21 Now let $g \in Z_i$ and $g' \in Z_j$ for some i, j and let m_t be an arbitrary ma-
 22 chine action. The sum of the factor values for this choice of partitions may be
 23 calculated offline as follows. If $i \neq j$ then
 24

$$25 \sum_{g \in Z_i} P(g_{t+1} = g' | g_t = g, m_t) = |Z_i| \theta_{m_t, 2} \quad (7)$$

26 and

$$27 \sum_{g' \in Z_j} P(g_{t+1} = g' | g_t = g, m_t) = |Z_j| \theta_{m_t, 2}. \quad (8)$$

28 If $i = j$ then

$$29 \sum_{g \in Z_i} P(g_{t+1} = g' | g_t = g, m_t) = (|Z_i| - 1) \theta_{m_t, 2} + \theta_{m_t, 1} \quad (9)$$

30 and

$$31 \sum_{g' \in Z_j} P(g_{t+1} = g' | g_t = g, m_t) = (|Z_j| - 1) \theta_{m_t, 2} + \theta_{m_t, 1}. \quad (10)$$

32 In all cases, the sums are independent of the particular g and g' chosen to
 33 represent the partitions. To simplify the notation, the factor value as a function
 34 of the partitions may now be defined by summing over the partition values.
 35 The partition factor value is then a simple multiple of the values calculated
 36 in Equations 7, 8, 9 and 10
 37
 38
 39
 40
 41
 42
 43
 44
 45
 46
 47
 48
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

$$f_a(Z_j, Z_i, m_t) = \sum_{g \in Z_i} \sum_{g' \in Z_j} f_a(g', g, m_t) \quad (11)$$

$$= \begin{cases} |Z_i|(|Z_i| - 1)\theta_{m_t,2} + |Z_i|\theta_{m_t,1} & \text{if } i = j \\ |Z_i||Z_j|\theta_{m_t,2} & \text{otherwise} \end{cases}. \quad (12)$$

As mentioned above, to simplify Equations 5 and 6 the messages into the factor must be constant before the factor values may be grouped together. Assume the messages are constant over the partition and denote the messages as a function of the partition - e.g. $\mu_{g_t \rightarrow f_a}(Z_i) = \mu_{g_t \rightarrow f_a}(g)$ for any $g \in Z_i$. The update equations can now be written as a function of the partition factor values:

$$\mu_{f_a \rightarrow g_{t+1}}(Z_j) = \sum_i \frac{1}{|Z_j|} f_a(Z_j, Z_i, m_t) \mu_{g_t \rightarrow f_a}(Z_i) \quad (13)$$

$$\mu_{f_a \rightarrow g_t}(Z_i) = \sum_j \frac{1}{|Z_i|} f_a(Z_j, Z_i, m_t) \mu_{g_{t+1} \rightarrow f_a}(Z_j). \quad (14)$$

These new equations can be used in the system by grouping together states when there is no evidence to distinguish them. The partitions are typically given by a collection of singletons along with a set representing the remaining options. When an observation arrives which distinguishes one of the states from the remaining options, that state is split off. Loopy belief propagation then proceeds as before, except that instead of multiplying by the normal factor value, the sums given above are used. The result is that a summation is only required for each partition rather than for every possible user goal.

An important requirement for these equations to remain valid is that the messages for each value must be equal. This means that a uniform prior must be used over the values within a partition. It may be possible to assume a non-uniform prior in certain special cases, as for example in the HIS system. However, an analysis of the exact circumstances where non-uniform priors may be used is left for future work.

2.4.2 General case

The use of partitions in loopy belief updates can be generalized to arbitrary factors. This allows the grouped LBP approach to be applied with validity nodes and also in the connections between the goal node and the user act. The details of the general algorithm are given in Appendix A.

2.5 Further efficiency gains

The grouped form of loopy belief propagation provides a significant reduction in calculation time. The algorithm still, however, requires an iteration over all combinations of the partitions. In certain special cases, this iteration can be simplified as well.

The case that is important for dialogue is the one discussed previously in Section 2.4.1. In this case, has the property that the factor sum is constant whenever the partitions, Z_i and Z_j , are different. Hence, this term can be factored-out from Equation 13 to give:

$$\begin{aligned}\mu_{f_a \rightarrow g_{t+1}}(Z_j) &= \frac{1}{|Z_j|} \sum_i \mu_{g_t \rightarrow f_a}(Z_i) f_a(Z_j, Z_i, m_t) \\ &= ((|Z_j| - 1)\theta_{m_t,2} + \theta_{m_t,1}) \mu_{g_t \rightarrow f_a}(Z_j) \\ &\quad + \theta_{m_t,2} \sum_{i \neq j} |Z_i| \mu_{g_t \rightarrow f_a}(Z_i).\end{aligned}$$

In loopy belief propagation, the effect of the messages on the estimated marginal distributions is independent of their scale and it is often necessary to rescale messages to avoid numerical issues. This can be exploited to further simplify the update equation. For grouped LBP, a suitable rescaling is to set:

$$\sum_i |Z_i| \mu_{g_t \rightarrow f_a}(Z_i) = 1. \quad (15)$$

Once the messages are scaled in this way, the update becomes:

$$\begin{aligned}\mu_{f_a \rightarrow g_{t+1}}(Z_j) &= ((|Z_j| - 1)\theta_{m_t,2} + \theta_{m_t,1}) \mu_{g_t \rightarrow f_a}(Z_j) \\ &\quad + \theta_{m_t,2} (1 - |Z_j| \mu_{g_t \rightarrow f_a}(Z_j)) \\ &= \theta_{m_t,1} \mu_{g_t \rightarrow f_a}(Z_j) + \theta_{m_t,2} (1 - \mu_{g_t \rightarrow f_a}(Z_j)).\end{aligned}$$

This update can therefore be done without any iteration through combinations of partitions. Instead, one only has to iterate through the set of partitions, resulting in a large reduction in computation. The update for the messages to the previous state can be optimised in a similar way.

2.6 Efficiency analysis

To formally analyse the efficiency of the above optimised form LBP, consider the updates for a factor with K different dependent variables, each having N values. Suppose further that the values for each node may be split into P indistinguishable partitions. The grouped form of loopy belief propagation must iterate through all combinations of the partitions and will therefore take $\mathcal{O}(P^K)$ calculations. Some extra partitioning machinery is required to ensure that sufficient partitions are included. This requires only a single pass through the current list of partitions so does not affect the complexity significantly. When the values are not grouped, the same factor updates would require $\mathcal{O}(N^K)$ calculations.

When a constant probability of change is assumed, there is no longer the need to go through all the combinations of the variable values. When values are grouped this results in a decrease by a factor of P . Goal updates in the model will therefore take $\mathcal{O}(P)$ calculations rather than $\mathcal{O}(P^2)$.

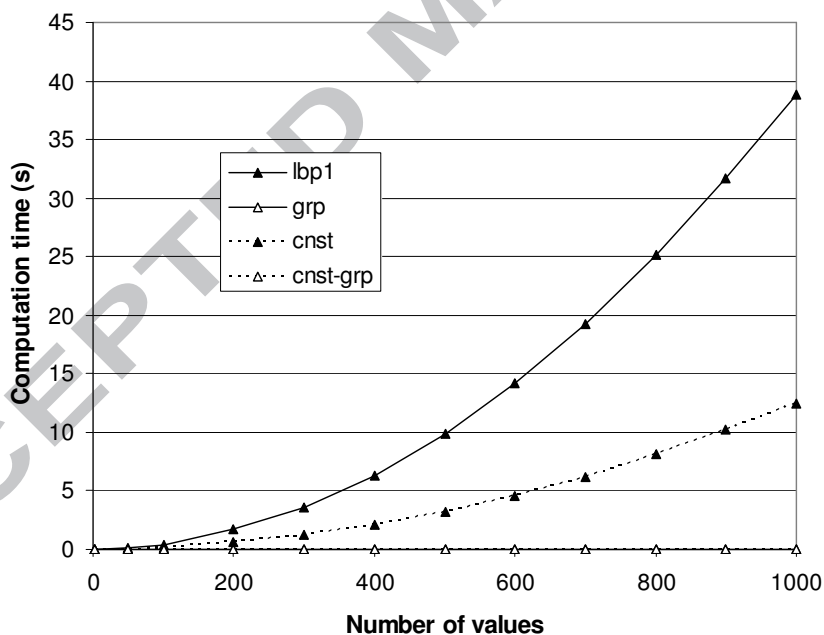


Fig. 5. Computation times of various algorithms for the marginal distribution after 10 time-slices in a Bayesian network with 1 node. Further details of the network are given in Appendix B.1. 500 samples were drawn from the network and the mean calculation time is plotted. The two algorithms using grouping take a maximum computation time of 0.04s and are therefore indistinguishable on the graph. The standard errors of all estimates are below 0.01s.

2.6.1 Experimental evaluation of efficiency

An experimental evaluation of the algorithms' efficiency was performed by measuring the computation times for several simple models defined in Appendix B.1⁵. Four algorithms were compared:

- The one time-slice LBP algorithm (lbp1). This is equivalent to the Boyen-Koller algorithm,
- The one time-slice LBP algorithm assuming a constant probability of change (cnst),
- The one time-slice LBP algorithm with grouping of values (grp),
- The one time-slice LBP algorithm assuming a constant probability of change and grouping values (cnst-grp).

The first experiment used a model with a single concept ($K = 1$) and the results are given in Figure 5. The computation times clearly follow the trends suggested by the theory above. As the number of values increases, the computation time of the standard algorithm increases with the square of the number of values. Assuming a constant probability of change reduces the computation time to a linear increase. The further reduction achieved by grouping values is dramatic since so few partitions are required that there is little computation to be done.

Figure 6 shows an evaluation of the proposed algorithms on models with more dependencies. In this figure, a tree structure of depth 2 and varying branching factors is used to represent the hidden goal. This gives an abstraction of the hierarchical dependencies often present in dialogue systems. Higher-level node values in this network can be grouped only rarely since there are too many other nodes which depend on them. One can see how this results in the assumption of a constant probability of goal changes being more effective than grouping. Grouping does, however, still improve performance.

These results show how the two techniques of grouping and assuming constant changes work together to improve efficiency. Grouping reduces the computation due to large numbers of values for a particular node. Assuming constant changes in the goal helps to reduce the computation required for nodes with more dependencies.

The factors connecting the sub-goals to the sub-acts can also use the grouped form of LBP by ensuring the partitions have equal user act probabilities. In most cases, this will not be a significant assumption. For example, it is reasonable to assume that the probabilities of the user saying "I would like Chinese food" given that they want Indian or Italian food would be the same

⁵ All experiments were performed on an Intel Core 2 Quad CPU, 2.4GHz processor with 4GB of RAM.

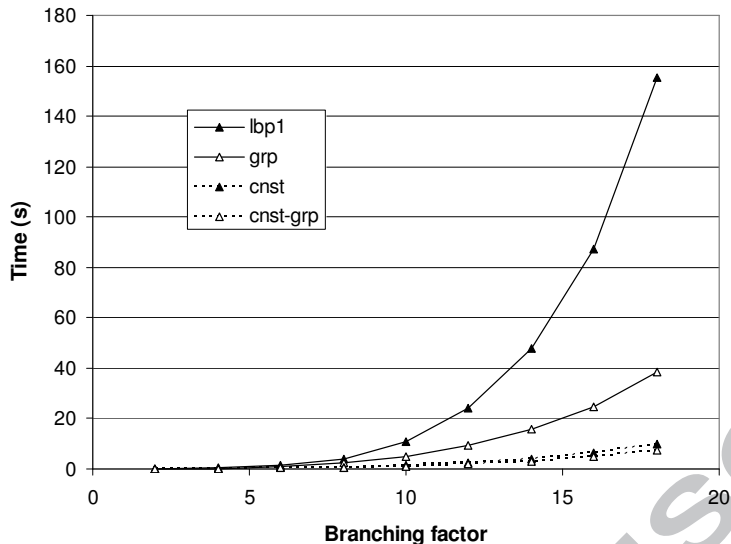


Fig. 6. Computation times of various algorithms for 10 time-slices of a Bayesian network. The network has a tree structure of depth 2 with varying branching factors. Further details of the network are given in Appendix B.1. 500 samples from each network were drawn and the mean calculation times are plotted. The standard errors of all estimates are below 0.7s. The standard errors for estimates of the *cnst* and *cnst-grp* algorithms are below 0.07s.

(i.e equally unlikely).

Once all possible simplifications have been included into a dialogue system, all of the important factor updates will require linear time in the number of concept-level partitions. The connection between the sub-acts and the overall act is deterministic so it requires only an iteration through the overall act values. The factors connecting history nodes have very few values to iterate through and do not contribute significantly to computation time. Further dependencies can also be modeled if necessary, but will bring associated increases in complexity.

2.6.2 Comparison with other approaches to belief state updating

When the unfactorised goal is used for belief updating then the number of partitions at the slot-level must be multiplied to obtain the number of unfactorised goal partitions. If the system state contains P partitions for each of K slots, the Hidden Information State model will require $\mathcal{O}(P^K)$ calculations to do the update.

$\mathcal{O}(P^K)$ is the same computation time as the proposed algorithms with all dependencies included. The difference is that belief updates in the HIS model are calculated exactly on the assumption that user goals do not change. The algorithms presented in this paper allow for changes in the user goal but

1 necessitate that approximations be made. If the entire goal is considered as a
 2 single node then the calculation times and accuracy are essentially the same.
 3 One can therefore think of the grouped belief propagation algorithm as an
 4 extension of the HIS approach to systems with conditionally independent slots
 5 and changes in the user goal.
 6

7 An alternative approach to updating the belief state is to use particle filters
 8 (Williams, 2007b). Comparing this method with the approach here is difficult
 9 since the computation needed for a particle filter depends crucially on the
 10 number of particles used. For a given level of accuracy, this is highly problem
 11 dependent and hence a fair comparison is left for future work.
 12
 13
 14

15 3 Policy Design

16
 17
 18
 19
 20 The essential role of the dialogue manager is to decide what action to take
 21 at each turn. This depends on the policy, which maps belief states, $b \in \mathcal{B}$,
 22 to probabilistic action choices. The structure of the policy will depend on the
 23 form of the belief state b .
 24

25
 26 In commercial systems, the policy is usually hand-crafted. At each point in
 27 the call flow, the dialogue designer chooses which action the system should
 28 take. Hand-crafted policies are also common for systems which use a statistical
 29 approach to uncertainty. After each turn the system updates its belief using
 30 probability theory, for example using the techniques described in Section 2.
 31 The policy then maps these updated beliefs to system actions. This policy can
 32 be hand-crafted, using decision trees based on selected elements of the belief
 33 state. Horvitz and Paek (1999) have suggested the Value of Information as
 34 a suitable quantity for basing these hand-crafted decisions. However, hand-
 35 crafting policies is time consuming and can result in sub-optimal decisions.
 36 Automatic policy learning is clearly preferable and this is the subject of the
 37 rest of this section.
 38
 39
 40
 41
 42
 43

44 3.1 Markov Decision Processes

45
 46
 47 The model traditionally used for policy learning is a Markov Decision Process
 48 (MDP), in which the belief state sequence, b_0, b_1, \dots, b_T , is finite and Markov⁶.
 49 The feature of an MDP which enables learning is the reward function, R ,
 50
 51

52 ⁶ The system is Markov when the probability of b_{n+1} given b_n is conditionally
 53 independent of b_i , for all $i < n$. The system is said to be *episodic* when the belief
 54 state sequence is finite. Non-episodic MDPs also require the definition of a discount
 55 rate, γ , which will be omitted here.
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

which defines the immediate reward for being in a particular state and taking a particular action, m . The aim of policy learning is then to optimise the expected future reward. This expected future reward is known as the *value function* which depends on the policy, π , and the start state, b :

$$V^\pi(b) = \sum_m \pi(b, m)R(b, m) + \sum_m \sum_{b'} \pi(b, m)P(b'|b, m)V^\pi(b'). \quad (16)$$

Several other quantities are useful when working with MDPs. The Q-function, $Q^\pi(b, m)$, is the expected future reward obtained by starting with a particular action and then following the policy. The advantage function, $A^\pi(b, m)$, is the difference between the Q-function and the value function. The occupancy frequency, $d^\pi(b)$, gives the expected number of times each state is visited⁷. These three quantities are given by the following equations:

$$Q^\pi(b, m) = R(b, m) + \sum_{b'} P(b'|b, m)V^\pi(b'), \quad (17)$$

$$A^\pi(b, m) = Q^\pi(b, m) - V^\pi(b), \quad (18)$$

$$d^\pi(b) = \sum_{t=0}^{\infty} P(b_t = b). \quad (19)$$

In the case of a POMDP, the reward is traditionally written as a function of the environment state s rather than the belief state b . The Markov assumption is also usually made on the environment states rather than the belief states. Fortunately, one can show that if the environment states are Markov then the belief states are too. A belief state reward may be derived from the environment state reward as follows:

$$R(b, m) = \sum_{s \in \mathcal{S}} b(s)R(s, m). \quad (20)$$

The belief state reward can be used to define a continuous state MDP that is closely related to the original POMDP. The internal state space of this MDP is identical to the belief state of the POMDP and policies will optimise the MDP only if they also optimise the original POMDP (Kaelbling et al., 1998). Throughout this paper, optimisation is performed on this belief state MDP instead of the original POMDP. This allows the techniques to be generalised for use in both MDP and POMDP frameworks, or even in hybrid approaches.

For the purpose of learning, it is useful to consider policies that are parameterised according to some parameter vector θ . The policy therefore determines

⁷ This is known to be finite because the system is episodic.

1 the probability of an action as a function of the parameters and the current
 2 belief state, $\pi(m|b, \theta)$. Policy learning can then be achieved by performing
 3 gradient descent on the policy whilst attempting to optimise $V^\pi(b)$.

4
 5 Before detailing the calculations, two techniques will be described for simpli-
 6 fying learning in a POMDP dialogue system. Firstly, the action set can be
 7 significantly reduced by encoding expert knowledge into summary actions, as
 8 shown in Section 3.2. Secondly, a parameterised policy can be used to factorise
 9 the effects of different concepts, as explained in Section 3.3.

13 3.2 Summary Actions

14
 15
 16
 17 In a dialogue system there are many cases where it is obvious to a sys-
 18 tem designer that one action is better than another. For example, it al-
 19 ways makes more sense to confirm the most likely value for a concept rather
 20 than any other value. Thus, if the probability of “food=Chinese” is greater
 21 than “food=Indian”, it would not make sense to attempt to confirm that
 22 “food=Indian”.
 23
 24

25
 26 Summary actions exploit this idea to significantly reduce the size of the action
 27 set. Along with the original machine action set, $\tilde{\mathcal{M}}$, a set of summary actions,
 28 \mathcal{M} , is defined. Note that \mathcal{M} is used here as the summary action set rather
 29 than the original set, as in Section 2, because it represents the set of actions
 30 which are actually used for learning. Each summary action, m , is defined by
 31 the system designer and represents a subset of $\tilde{\mathcal{M}}$. Given a summary action,
 32 m , and belief state, \mathbf{b} , a mapping back to the original action set, $F(m, \mathbf{b})$,
 33 is defined.⁸ In the example, “confirm food” action would be the associated
 34 summary act and given a belief state, the corresponding machine act in $\tilde{\mathcal{M}}$
 35 would be to confirm the most likely food type.
 36
 37
 38

39
 40 A new continuous state MDP can be defined using this summarised action
 41 set. The belief states are exactly as before but the actions are restricted to
 42 \mathcal{M} . Given a summary action and belief state, m and \mathbf{b} , the system simply
 43 takes the corresponding original action, $F(m, \mathbf{b})$. Rewards and transitions are
 44 defined as before using this corresponding action. The system is still an MDP,
 45 but the action set is reduced.
 46
 47

48
 49 This process of using summary actions allows the system designer to encode
 50 expert knowledge to allow learning to proceed more quickly. The optimal pol-
 51 icy for this summarised version of the MDP will be close to the optimal policy
 52 for the full MDP provided that the inverse mapping, F , is reasonably accurate.
 53

54
 55 ⁸ The belief state is now shown in bold to indicate its concrete representation as a
 56 vector of probability values.
 57
 58

1 The use of summary actions is based on the summary POMDP idea proposed
 2 by Williams and Young (2005). The difference is that the summary POMDP
 3 factors both the state and actions according to a set of slots. This necessitates
 4 a calculation of the state dynamics at a slot level and makes learning between
 5 actions that affect different slots impossible. Also, the changed dynamics of
 6 the summary POMDP require special learning algorithms. In contrast, the ap-
 7 proach here does not change the state, makes no change to the state dynamics
 8 and preserves the properties of a standard MDP so no special techniques are
 9 required.
 10

11 3.3 Component-based policies for dialogue

12
 13
 14 When automating decision making, it is useful to be able to separate out the
 15 effects of different components. Using parameterised policies provides a simple
 16 method of implementing this.
 17

18 A standard parameterisation for the policy is to use a softmax function. For
 19 each summary machine action, m , a *basis function* is defined, ϕ_m , which is
 20 a vector function of the belief state, \mathbf{b} . Each basis function identifies a set
 21 of features from the belief state that are important for decision making and
 22 its value is computed as an arbitrary function of these features. Along with
 23 the parameters, θ , these basis functions determine the likelihood of taking an
 24 action according to the following equation:
 25

$$26 \pi(m|\mathbf{b}, \theta) = \frac{e^{\theta \cdot \phi_m(\mathbf{b})}}{\sum_{m'} e^{\theta \cdot \phi_{m'}(\mathbf{b})}}. \quad (21)$$

27 This parameterisation provides a simple method of separating the effects of
 28 each component on the overall probability of taking each action. For a dialogue
 29 system, each concept, $c \in \mathcal{C}$, can be associated with a concept-level basis
 30 function, $\phi_{m,c}(\mathbf{b})$ and any additional global information can be encoded in an
 31 overall basis function, $\phi_{m,*}(\mathbf{b})$. The full basis function can then be separated
 32 into components as follows:
 33

$$34 \phi_m(\mathbf{b})^\top = [\phi_{m,1}(\mathbf{b})^\top, \dots, \phi_{m,K}(\mathbf{b})^\top, \phi_{m,*}(\mathbf{b})^\top]. \quad (22)$$

35 Each concept-level function must encode a set of features that depends on
 36 the machine action. As an example, a typical ‘‘slot-filling’’ dialogue might
 37 separate $\phi_{m,c}(\mathbf{b})$ into components which each depend on how the machine
 38 action affects the given slot, c . A straightforward way of achieving this is
 39 to use an indicator function along with a set of sub-components. For ex-
 40

ample, if a set of N_c different sub-components are defined then $\phi_{m,c}(\mathbf{b})^\top = [\phi_{m,c,1}(\mathbf{b})^\top, \dots, \phi_{m,c,N_c}(\mathbf{b})^\top]$, where all of these *sub basis functions*, $\phi_{m,c,i}(\mathbf{b})$, are $\mathbf{0}$ except for the sub-component corresponding to the given machine action. This nonzero component will be a function of the beliefs, decided by the system designer. The component should give an indication of the effect of the beliefs on action decisions for that particular slot.

As an example, consider a system where the only available actions are to request the value of different slots. The system designer may decide that the only important information for policy making is the probability of the most likely value for each concept, which is denoted l_c . This probability is to be multiplied by one parameter, $\theta_{c,1}$, if the action is to request this slot and a different parameter, $\theta_{c,2}$, if the action is to request a different slot. A suitable set of functions is then:

$$\phi_{m,c}(\mathbf{b}) = \begin{cases} \begin{pmatrix} l_c \\ 0 \end{pmatrix} & \text{if } m \text{ requests slot } c \\ \begin{pmatrix} 0 \\ l_c \end{pmatrix} & \text{otherwise} \end{cases} . \quad (23)$$

In practice, a grid-based approach may well be preferred for these sub basis functions and this is the approach used in the Tourist Information System described later in Section 4. In the grid-based approach, a value of 1 is assigned to the component of the grid given by the beliefs and a value of 0 is assigned to all other components. For example, the probabilities of the two most likely slot values might be formed into a tuple with tuples grouped together whenever the most likely slot value is less than 0.5 and otherwise identified with the closest point in the set $\{(1.0, 0.0), (0.8, 0.0), (0.8, 0.2), (0.6, 0.0), (0.6, 0.2), (0.6, 0.4)\}$. A given sub basis function, $\phi_{m,c,i}(\mathbf{b})$, will then map to a 7-dimensional vector with 1 assigned to the dimension associated with the corresponding group and 0 everywhere else. Other alternatives are possible, for example using two components: the entropy of the slot and the probability of the most likely option could be used.

The final basis function that must be defined is the overall basis function $\phi_{m,*}(\mathbf{b})$. This function allows the system designer to incorporate global knowledge about the overall state in the decision making. Again taking as an example the system described in Section 4, a useful global metric is the number of database items that match a given request given the current most likely user goal. This number is determined and grouped into one of “no venues”, “one venue”, “a few venues” or “many venues”. $\phi_{m,*}(\mathbf{b})$ is then a 4-dimensional vector with 1 in the dimension associated with the relevant group and 0 otherwise. Policy learning is used to decide how important this information is.

Other knowledge about the overall state of the system could be incorporated in a similar way.

3.4 Natural Actor Critic

Once the policy has been parameterised using a suitable structure, the system must learn the parameters that maximise the expected future reward. There are various possible approaches to performing this optimisation but one that has worked well in the framework presented here is the Natural Actor Critic (NAC) algorithm (Peters et al., 2005), which performs policy optimisation using a modified form of gradient descent.

Traditional gradient descent iteratively adds a multiple of the gradient to the parameters being estimated. This is not necessarily the most effective gradient to use because the Euclidean metric may not be appropriate as a measure of distance. In general, the parameter space is described as a Riemann space, and a metric tensor, G_{θ} , can be defined such that for small changes in the parameters θ , the distance is $|d\theta|^2 = d\theta^{\top} G_{\theta} d\theta$. The traditional Euclidean distance is then given by a metric tensor equal to the identity matrix. In optimising an arbitrary loss function, $L(\theta)$, Amari (1998) shows that for a general Riemann space the direction of steepest descent is in fact $G_{\theta}^{-1} \nabla_{\theta} L(\theta)$. This gradient is known as the *natural gradient* to contrast it with the *vanilla gradient* traditionally used.

The optimal metric tensor to use in a statistical model is typically the Fisher Information Matrix which has been shown to give distances that are invariant to the scale of the parameters (Amari, 1998). Given a probability distribution $p(x|\theta)$, the Fisher Information is the matrix G_{θ} such that $(G_{\theta})_{ij} = \mathbb{E}(\frac{\partial \log p(x|\theta)}{\partial \theta_i} \frac{\partial \log p(x|\theta)}{\partial \theta_j})$. Peters et al. (2005) show that the Fisher Information matrix for an MDP model is:

$$G_{\theta} = \int_{\mathcal{B}} d^{\pi}(b) \int_{\mathcal{M}} \pi(m|b, \theta) \nabla_{\theta} \log \pi(m|b, \theta) \nabla_{\theta} \log \pi(m|b, \theta)^{\top} dm db. \quad (24)$$

The direction of steepest descent is the inverse of this matrix multiplied by the traditional vanilla gradient. The vanilla gradient for an MDP is given by the Policy Gradient Theorem of Sutton et al. (2000):

$$\nabla_{\theta} V(b_0) = \int_{\mathcal{B}} d^{\pi}(b) \int_{\mathcal{M}} A^{\pi}(b, m) \pi(m|b, \theta) \nabla_{\theta} \log \pi(m|b, \theta) dm db. \quad (25)$$

Equation 25 depends crucially on two quantities which are difficult to compute: the advantage function, $A^\pi(b, m)$ defined in Equation 18, and the occupancy frequency, $d^\pi(b)$ defined in Equation 19.⁹ The process of integrating over the occupancy frequency may be approximated with sampling techniques as explained in Section 3.4.1. At the same time, the advantage function is replaced by a compatible approximation. The resulting equations for the natural gradient then simplify dramatically since the Fisher Information matrix cancels. More specifically, Peters et al. (2005) shows that if an approximate advantage function, $a_{\mathbf{w}}(b, m)$ is chosen such that:

$$a_{\mathbf{w}}(b, m) = \nabla_{\boldsymbol{\theta}} \log \pi(m|b, \boldsymbol{\theta}) \cdot \mathbf{w}, \quad (26)$$

where \mathbf{w} minimises the average squared approximation error, i.e.

$$\frac{\partial}{\partial \mathbf{w}} \int_{\mathcal{B}} d^\pi(b) \int_{\mathcal{M}} \pi(m|b, \boldsymbol{\theta}) (A^\pi(b, m) - a_{\mathbf{w}}(b, m))^2 dm db = \mathbf{0}, \quad (27)$$

then the required natural gradient is given by

$$G_{\boldsymbol{\theta}}^{-1} \nabla_{\boldsymbol{\theta}} V(b_0) = \mathbf{w}. \quad (28)$$

This gradient can be used to develop an algorithm with relatively strong convergence properties. The system iterates between evaluating the policy and improving it. The evaluation, called the *critic* step, must estimate the approximate advantage function. Once a suitable function is estimated an *actor improvement* is made by changing the parameters by a multiple of the natural gradient. The algorithm, called the Natural Actor Critic (NAC), is guaranteed to converge to a local maximum of the value function provided that Equation 27 is satisfied at every step and certain other minor conditions are met.

3.4.1 Sampling methods

Equation 27 requires a minimisation which does not generally have an analytical solution. As a result, the NAC algorithm uses sampling to *estimate* the gradient. Assuming that the dialogues are numbered $n = 1 \dots N$ and the turns are numbered $0, \dots, T_n - 1$, the equivalent sampled equation to minimize is:

$$\sum_{n=1}^N \sum_{t=0}^{T_n-1} |A(b_{n,t}, m_{n,t}) - a_{\mathbf{w}}(b_{n,t}, m_{n,t})|^2. \quad (29)$$

Since the system does not actually have estimates for the true advantage function, $A(b, m)$, the rewards in a dialogue must be grouped together in

⁹ The occupancy frequency is also sometimes called the *state distribution* (Peters et al., 2005).

order to obtain suitable estimates. The sum of the rewards gives an unbiased estimate of the sum of the advantages and the initial value function. The resulting equation is as follows:

$$\sum_{t=0}^{T_n-1} A(b_{n,t}, m_{n,t}) + V^\pi(b_0) \approx \sum_{t=0}^{T_n-1} r(b_{n,t}, m_{n,t}). \quad (30)$$

The task of solving for a suitable approximation is now changed from minimising Equation 29 to minimising the average error per dialogue. If J represents an estimate for the value function of the start state then the error is:

$$\sum_{n=1}^N \left| \sum_{t=0}^{T_n-1} A(b_{n,t}, m_{n,t}) - \sum_{t=0}^{T_n-1} a_{\mathbf{w}}(b_{n,t}, m_{n,t}) \right|^2 \quad (31)$$

$$= \sum_{n=1}^N \left| \sum_{t=0}^{T_n-1} r(b_{n,t}, m_{n,t}) - \left(\sum_{t=0}^{T_n-1} \nabla_{\boldsymbol{\theta}} \log \pi(m_{n,t} | b_{n,t}, \boldsymbol{\theta}) \cdot \mathbf{w} - J \right) \right|^2. \quad (32)$$

This is a simple least squares regression problem and can be solved using standard techniques. When the policy parameters are changed in an actor improvement step, the previous dialogues no longer give valid estimates for the approximation. In order to reduce this effect, previous dialogues are successively deweighted. Algorithm 2 gives a complete implementation of episodic NAC.

Algorithm 2 Episodic Natural Actor Critic

for each dialogue, n **do**

 Execute the dialogue according to the current policy π .

 Obtain the sequence of states, $b_{n,t}$, and machine actions, $m_{n,t}$.

 Compute statistics for the dialogue:

$$\boldsymbol{\psi}_n = \left[\sum_{t=0}^{T_n} \nabla_{\boldsymbol{\theta}} \log \pi(m_{n,t} | b_{n,t}, \boldsymbol{\theta})^\top, 1 \right]^\top,$$

$$R_n = \sum_{t=0}^{T_n} r(b_{n,t}, m_{n,t}).$$

Critic Evaluation: Choose \mathbf{w} to minimize $\sum_n (\boldsymbol{\psi}_n^\top \mathbf{w} - R_n)^2$.

Actor-Update: If \mathbf{w} has converged:

 Update the policy parameters: $\boldsymbol{\theta}_{n+1} = \boldsymbol{\theta}_n + \mathbf{w}_0$, where $\mathbf{w}^\top = [\mathbf{w}_0^\top, J]$.

 Deweight all previous dialogues: $R_i \leftarrow \gamma R_i$, $\boldsymbol{\psi}_i \leftarrow \gamma \boldsymbol{\psi}_i$ for all $i \leq n$.

end for

3.5 Comparison with other algorithms

The natural actor critic is one of several algorithms that could be used to learn the policy. Other popular learning algorithms include TD-learning, SARSA-learning and various improvements of point-based value iteration (Sutton and Barto, 1998; Shani et al., 2008; Williams and Young, 2006a).

One particularly important characteristic of the NAC algorithm is that it is *model-free* rather than *model-based*. This means that the algorithm learns from sample dialogues rather than using the model inherent in the Bayesian network. The major advantage of this is that it is much easier to build a good simulator of the environment than it is to build a good probabilistic model (Sutton and Barto, 1998). Another important reason for using a model-free approach is that the observations in a dialogue system are not finite. The use of confidence scores requires a continuous set of observations, unless they are discretised. This is the main reason for using NAC over model-based techniques such as PBVI.

NAC's strong convergence properties make it preferable to other model-free approaches, such as TD- and SARSA learning. The use of gradient descent results in small changes to the policy. This allows for a much more gradual learning process, which reduces the chances of reaching a local optimum (Peters et al., 2005).

4 Evaluation

4.1 An example task

The framework described above was evaluated by comparing it with a conventional MDP system for the tourist information domain, in which users may ask about hotels, restaurants, bars and amenities in a fictitious town. The dialogue is mixed-initiative, meaning that the system may ask specific questions while users may provide information that was not requested or ask their own questions. Four example systems were built in order to compare differences between hand-crafted and learned policies as well as the effects of using a statistical approach to model uncertainty (i.e. the four frameworks of Figure 1).

The user may speak about nine concepts in their attempts to find a suitable venue. These are: name of the venue, type of venue, area, price range, nearness to a particular location, type of drinks, food type (for restaurants), number

1 of stars (for hotels and restaurants) and music (for restaurants and bars).
2 Once a suitable venue is found the user may ask about four further concepts:
3 address, telephone number, a comment on the venue and the price (for hotels
4 and restaurants). The database contains 47 different venues.
5

6 A simulated user was built for the task using the agenda-based approach
7 of Schatzmann et al. (2007). The simulated user receives the semantic level
8 output of the dialogue system and gives a response deemed to reflect the way
9 a human user would behave. The simulated user is a useful tool for developing
10 dialogue systems as it supports both policy learning and evaluation. The same
11 simulator was used for both learning and for evaluation.
12
13

14
15 Outputs of the simulated user were passed through a simulated error channel
16 before being input to the dialogue manager. The error channel makes use of
17 a *confusion rate*, ρ , which determines the amount of error produced. At each
18 turn, three alternatives are produced which are each correct with probability
19 $1 - \rho$. These alternatives are then collated into an N-best list of user acts
20 and confidence scores are attached. Most experiments used a confusion rate
21 of 40%. A typical run of 1000 dialogues at this confusion rate results in a top
22 accuracy of 80.9, an oracle accuracy of 96.3 and an ICE score of 0.96. The
23 top accuracy gives a measure of the usefulness of the most likely hypothesis,
24 the oracle accuracy gives the accuracy for the best hypothesis in the N-best
25 list and the ICE score measures the overall usefulness of the confidence scores.
26 Further details on these metrics may be found in Thomson et al. (2008).
27
28
29
30

31 Evaluation was based on four metrics:
32
33

- 34 • **Reward:** The reward is 20 less the number of dialogue turns for a successful
35 dialogue and 0 less the number of turns for an unsuccessful one. A dialogue
36 was considered successful if a suitable venue was offered and all further
37 pieces of information were given. In the case where no venue matched the
38 constraints, the dialogue was deemed successful if the system told the user
39 that no venue matched and a suitable alternative was offered.
40
- 41 • **Objective Success:** A dialogue was objectively successful if a suitable
42 venue was offered and all further pieces of information were given. If no
43 venue matched the constraints then the dialogue was successful if the user
44 was told that there was no matching venue.
45
- 46 • **Objective Score:** If the dialogue was objectively successful, then the score
47 was 20 less the number of turns until success. Otherwise, it was 0 less the
48 number of dialogue turns.
49
- 50 • **Subjective Success:** For testing with real users only, a dialogue was deemed
51 subjectively successful if the user believed that their requirements had been
52 satisfied.
53
54
55

56 The difference in the definitions of success for the reward and the objective
57
58
59

score is due to a lack of information when testing with human users. From the dialogue transcriptions, it is not possible to tell how a human user changes their goals when told there is no matching venue whereas the simulated user provides this information explicitly.

4.2 The four systems

A Bayesian network for the tourist information domain was built using the framework for Bayesian updates of dialogue state (BUDS) described in Section 2.2. A user goal node, a user act node and a history node were implemented for each of the nine concepts used to find suitable venues. User act nodes and history nodes were added for the four extra concepts that could be asked. An extra “method” node was also included to capture the approach being used by the user to find the venue. This enables a venue to be found by name, according to a set of constraints or for a set of alternative venues to be requested given some constraints. Figure 7 shows a graphical representation of the concepts in the system and their dependencies.

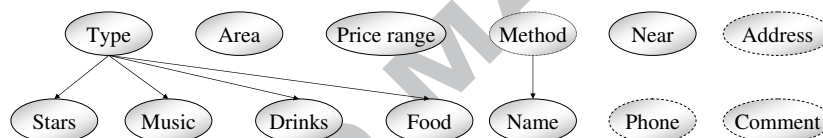


Fig. 7. A graphical representation of the concepts for one time slice in the BUDS experimental system. The “phone”, “address” and “comment” concepts are ignored at the goal-level and only maintain nodes at the act- and history-levels. The other concepts have corresponding nodes at the goal level with a Bayesian Network structure as depicted in the diagram.

Policy learning for the POMDP system was implemented using the episodic Natural Actor Critic algorithm and a component-based policy, as explained in Section 3. One summary action enabled the system to offer information about a venue while the others asked the user to confirm, request or select one of two options for each of the different concepts.

Although the NAC algorithm can be performed online, many thousands of dialogues are required hence online learning with real humans is not practical.¹⁰ Instead, the system interacted with the agenda-based simulator discussed above. The same reward function was used as for evaluation, thus a 20 point reward was given for full completion of a dialogue and 1 point was subtracted for each dialogue turn. Figure 8 shows how the average reward

¹⁰ Although human-based on-line learning at these call volumes might be practical for some of the providers of large commercial spoken dialogue systems!

improves as more dialogues are completed. After 250,000 dialogues the policy was fixed to yield the trained policy denoted as BUDS-TRA below ¹¹.

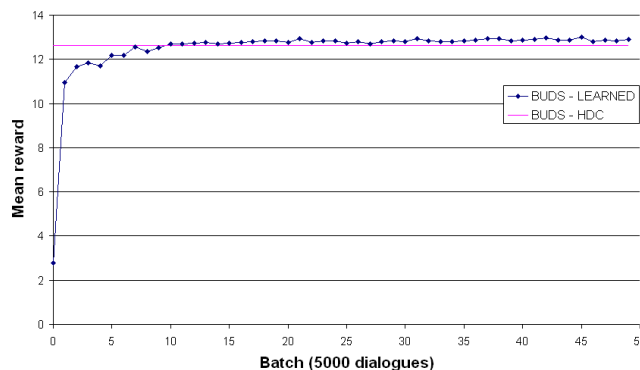


Fig. 8. Average reward over 5000 dialogues during training of the policy for use with the BUDS framework. The horizontal line gives the reward obtained by the hand-crafted policy with the same model for updating the belief state (BUDS-HDC).

A hand-crafted policy for the BUDS framework was also implemented by building up overall actions from sub-actions for each concept (labeled BUDS-HDC). All slots where the most likely value had probability greater than 0.8 were accepted. When the accepted slots gave sufficient information for a venue recommendation, the system gave the information. Otherwise the system would implicitly confirm, explicitly confirm, ask the user to select between two options or request (in that order) one of the low probability slots. In cases where the most likely value had probability greater than 0.5, the slot was confirmed or implicitly confirmed. When the probability was less than 0.4 it was requested. Otherwise the user was asked to select between the two most likely options.

Comparisons were made with two MDP systems which had a belief state composed of the same concepts as used in the BUDS-based systems. Each concept or slot could take three state values: unknown, known or confirmed. State transitions were based only on the most likely user act since it is impractical to use a full N-best list of user acts in an MDP system. No information from the confidence scores was maintained. A hand-crafted policy was built to either offer information or request or confirm a slot based on this state (this system is labeled HDC). A policy using the same belief state was also trained using standard reinforcement learning techniques (labeled MDP). These MDP systems were built by another researcher in another project and considerable effort had been expended to make them as competitive as possible (Schatzmann, 2008).

¹¹ The error model for training used a fixed confusion rate of 0.4 whereas the tests were conducted over a range of error rates. Thus, the test and training conditions were unmatched.

4.3 Simulation results

Figure 9 shows the results of a simulated comparison of the four systems. In addition, a variant of the BUDS based system with the hand-crafted policy was also tested (BUDS-HDC-TP1) which used only the most likely output from the simulated error channel. This allowed the benefit of using the full N-best semantic output to be assessed.

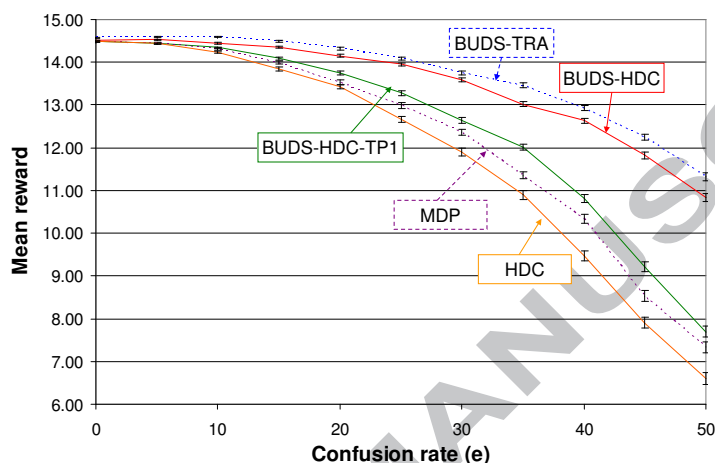


Fig. 9. Simulated comparison of the four frameworks for dialogue management. Each point gives the mean reward for 5000 simulated dialogue. Error bars show one standard error on each side of the mean.

The figure highlights some important features of the different policies. At low error rates all policies perform equally well, but at higher error rates a consistent ordering of the approaches emerges. Firstly, the two systems that make use of the full N-best list of user acts show significantly better performance. The importance of the N-best list is particularly evident in the comparison between the BUDS-HDC and BUDS-HDC-TP1 systems, which use all the same components except for the size of the N-best list. Secondly, the use of a statistical model of uncertainty outperforms hand-crafted models of the belief state, even when only the most likely user act is used. By using a systematic approach, the system is better able to handle conflicting evidence from different dialogue turns. Finally, the learned policies give significant improvements over the hand-crafted policies indicating that the learning has been effective.

4.4 User trial

Various authors have commented that evaluations using simulated users are not necessarily good predictors of performance with real human users. The systems were therefore also evaluated in a user trial by 36 native English speakers,

1 none of whom had been involved in previous experiments with dialogue sys-
2 tems. Users were asked to find a venue according to a set of constraints and
3 then obtain some extra pieces of information. 48 different task scenarios were
4 used, including cases where there was no suitable venue, cases where there was
5 more than one venue and cases where there was exactly one matching venue.
6 For each system, every user was given one task from each of these three groups,
7 resulting in a total of 108 dialogues recorded for each system.
8

9
10 The robustness of the systems was evaluated by adding synthetic noise to the
11 speech signal before speech recognition. Three levels of noise were used: low,
12 medium and high corresponding to signal to noise ratios (SNR) of 35.3dB,
13 10.2dB and 3.3dB respectively. Each user was tested under all three noise
14 conditions for each system.
15

16
17 All four systems use the same speech recogniser, semantic decoder, output
18 generator and text-to-speech engine. Speech recognition is implemented using
19 the ATK toolkit with a tri-phone acoustic model and a dictionary of around
20 2000 in-domain words. The output of the recogniser is a 10-best list along with
21 sentence-level inference evidence scores, which are the sum of the sentence arc
22 log-likelihoods in the confusion network. A hand-crafted Phoenix-based parser
23 is used for semantic decoding. The confidence for each resulting dialogue act
24 is calculated by exponentiating the inference evidence, adding the score for
25 sentences resulting in the same dialogue act and renormalising so that the
26 sum is one. Simple template rules are used for output generation and text-to-
27 speech uses a diphone synthesis engine.
28
29
30
31

32 Overall results from the trial are given in Table 1. It is clear from the results
33 that the systems which use a statistical approach to modeling uncertainty give
34 a significant improvement in performance. Interestingly, the BUDS trained
35 policy outperformed the hand-crafted policy in terms of subjective success but
36 was worse on the objective metrics. An investigation of the transcripts showed
37 that the trained policy sometimes offered venues before all the information had
38 been given. While the simulated user would always give further information
39 to ensure its constraints were met, the human users did not necessarily do
40 this. In these cases, the human user often thought the dialogue was successful
41 when in fact some constraint was not met.
42
43
44

45
46 An analysis of the systems' robustness to noise was evaluated by separating
47 out the results for different noise levels ¹², as shown in Figure 10. The overall
48

49
50 ¹²The questions asked by the dialogue manager can influence both the Word Er-
51 ror Rate (WER) as well as the Semantic Error Rate (SER). It is therefore more
52 important to compare success under equivalent noise conditions than to compare
53 success under equivalent error rates. Further experiments, not included in this pa-
54 per, showed similar trends when the success rate was plotted as a function of the
55 error rate instead of the noise level.
56
57
58
59
60
61
62
63
64
65

System	OSR	OS	SSR
BUDS-HDC	0.84 ± 0.04	11.83 ± 0.94	0.84 ± 0.04
BUDS-TRA	0.75 ± 0.04	8.89 ± 1.09	0.88 ± 0.03
MDP-HDC	0.65 ± 0.05	7.10 ± 1.21	0.78 ± 0.04
MDP-TRA	0.66 ± 0.05	6.97 ± 1.23	0.81 ± 0.04

Table 1

Objective Success Rate (OSR), Objective Score (OS) and Subjective Success Rate (SSR) for the different systems. Error values give one standard error of the mean. They are calculated by assuming that success rates and rewards follow binomial and Gaussian distributions respectively.

trend is that the BUDS systems do outperform the traditional approaches but these differences are not significant due to the small number of samples for each noise level. The graph shows a large variability in success rates which makes the results difficult to interpret. The objective success rates gave similar trends, although the BUDS-HDC policy performed better than the BUDS-TRA policy.

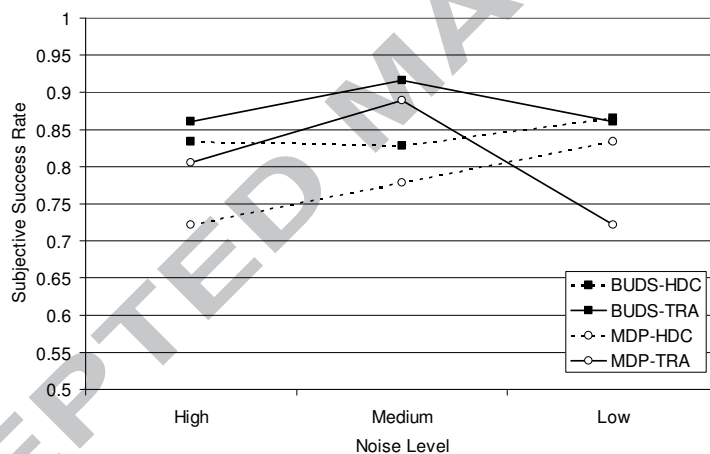


Fig. 10. Subjective success rates for the different systems at varying noise levels.

5 Conclusion

This paper has introduced a new POMDP-based framework for building spoken dialogue systems based on using Bayesian updates of the dialogue state (BUDS). The BUDS model gives an efficient, statistically motivated approach to modeling the internal system state. Various modifications of standard Bayesian network algorithms were presented which give sufficient efficiency improvements over the standard loopy belief propagation algorithm to enable a real-world dialogue system to operate efficiently in real-time. The belief updating

1 equations were also shown to be more efficient than alternative updating al-
2 gorithms used in POMDP-based dialogue systems.

3
4 Dialogue systems using the BUDS framework can use either learned or hand-
5 crafted policies. A learning algorithm based on the Natural Actor Critic al-
6 gorithm was presented as a feasible approach to policy optimisation. The fac-
7 torisation of the state space inherent in the BUDS model leads to a simple
8 separation of the policy parameters into components. Domain knowledge may
9 be incorporated by using summary actions to group together different actions
10 and increase the tractability of the policy learning.

11
12
13 The BUDS framework was evaluated in both simulations and in a user trial.
14 The two systems using the BUDS approach outperformed both the MDP and
15 finite state systems. This difference was significant in both simulations and
16 the user trial. The effect of policy learning with BUDS was not conclusive. In
17 simulations the learned policy clearly outperforms the hand-crafted one but
18 this did not translate into a conclusive improvement in the objective perfor-
19 mance with human users. Users did, however, *perceive* the learned policy to
20 be more successful than the hand-crafted policy.

21
22
23 One possible reason for the lower performance of the learned policy is the
24 difference between the training and testing environments. The simulator, error
25 channel and reward are all different between learning with the simulator and
26 testing with human users. Thus, over-tuning on the user simulator is highly
27 likely. Also, the policy is optimised for a reward which reflects the user's
28 opinion of whether the dialogue was successful. This is arguably closer to
29 the subjective success rather than to the objective success. Future learned
30 policies should be able to conclusively outperform hand-crafted policies by
31 using more appropriate rewards, user simulations, summary actions and basis
32 function features.

33
34
35 Another possible reason for the reduced performance with human users is the
36 choice of model parameters which were here estimated very roughly by inspect-
37 ing their effects on dialogue evolution with the user simulator. It would clearly
38 be preferable to learn these parameters from data by extending the Bayesian
39 network to include the parameters. Unfortunately, Loopy Belief Propagation,
40 would no longer be applicable because the parameter nodes are continuous
41 and so a more general algorithm, such as Expectation Propagation, would be
42 required. The authors are currently pursuing this topic.

43
44
45 It is clear from the paper that the use of a statistical approach to uncer-
46 tainty can yield significant improvements in performance. Policy learning can
47 also achieve competitive policies with little human intervention. The BUDS
48 framework provides an approach that is efficient enough to scale to real-world
49 problems, and which can outperform standard alternatives. Further work will
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 attempt to learn the model parameters from data and to deploy the system
2 in a live telephone application, showing that the benefits discussed here do
3 translate into real-world improvements.
4
5
6

7 Acknowledgements

8
9

10 This research was funded by a St John's Benefactors Scholarship, the UK EP-
11 SRC under grant agreement EP/F013930/1 and by the EU FP7 Programme
12 under grant agreement 216594 (CLASSIC project: www.classic-project.org).
13 The authors would like to thank Jost Schatzmann, who built the user simula-
14 tor and MDP based systems used here, as well as Milica Gašić, Filip Jurčićek,
15 Simon Keizer, Fabrice Lefèvre, François Mairesse, Ulrich Paquet and Kai Yu
16 for their help in the user trials and for useful comments and discussions.
17
18
19
20
21

22 References

23

- 24
25 Amari, S., 1998. Natural gradient works efficiently in learning. *Neural Com-*
26 *putation* 10 (2), 251–276.
27
28 Bishop, C., 2006. *Pattern Recognition and Machine Learning*. Springer.
29
30 Boyen, X., Koller, D., 1998. Tractable inference for complex stochastic pro-
31 cesses. In: *Proceedings of the 14th Conference on Uncertainty in Artificial*
32 *Intelligence (UAI)*. San Francisco: Morgan Kaufmann, pp. 33–42.
33
34 Bui, T., Poel, M., Nijholt, A., Zwiers, J., 2007. A tractable DDN-POMDP ap-
35 proach to affective dialogue modeling for general probabilistic frame-based
36 dialogue systems. In: Traum, D., Alexandersson, J., Jonsson, A., (eds), I. Z.
37 (Eds.), *Workshop on Knowledge and Reasoning in Practical Dialog Sys-*
38 *tems, International Joint Conference on Artificial Intelligence (IJCAI)*. Hy-
39 *derabad, India*, pp. 34–37.
40
41 Heskes, T., 2003. Stable fixed points of loopy belief propagation are minima
42 of the Bethe free energy. In: Becker, S., Thrun, S., Obermayer, K. (Eds.),
43 *Advances in Neural Information Processing Systems 15*. MIT Press, Cam-
44 *bridge*, pp. 359–366.
45
46 Horvitz, E., Paek, T., June 1999. A computational architecture for conver-
47 sation. In: *Proceedings of the Seventh International Conference on User*
48 *Modeling*. New York: Springer Wien, Banff, Canada, pp. 201–210.
49
50 Kaelbling, L. P., Littman, M. L., Cassandra, A. R., 1998. Planning and acting
51 in partially observable stochastic domains. *Artificial Intelligence* 101, 99–
52 134.
53
54 Kschischang, F., Frey, B., Loeliger, H., 2001. Factor graphs and the sum-
55 product algorithm. *IEEE Transactions on Information Theory* 47, 498–519.
56
57 Lemon, O., Georgila, K., Henderson, J., Stuttle, M., 2006. An ISU dialogue
58
59
60
61
62
63
64
65

- 1 system exhibiting reinforcement learning of dialogue policies: generic slot-
2 filling in the TALK in-car system. In: Proceedings of the European chapter
3 of the Association for Computational Linguistics (EACL).
- 4 Levin, E., Pieraccini, R., Eckert, W., 2000. A Stochastic Model of Human-
5 Machine Interaction for Learning Dialog Strategies. *IEEE Transactions on*
6 *Speech and Audio Processing* 8 (1), 11–23.
- 7 Meng, H., Wai, C., Pieraccini, R., November 2003. The use of belief networks
8 for mixed-initiative dialog modeling. *IEEE Transactions on Speech and Au-*
9 *dio Processing* 11 (6), 757–773.
- 10 Minka, T., 2001. A family of algorithms for approximate bayesian inference.
11 Ph.D. thesis, MIT.
- 12 Murphy, K., July 2002. Dynamic bayesian networks: Representation, inference
13 and learning. Ph.D. thesis, UC Berkeley, Computer Science Division.
- 14 Peters, J., Vijayakumar, S., Schaal, S., 2005. Natural actor-critic. In: *European*
15 *Conference on Machine Learning (ECML)*. Springer, pp. 280–291.
- 16 Pieraccini, R., Huerta, J. M., 2008. Where do we go from here? In: Dybkjr,
17 L., Minker, W. (Eds.), *Recent Trends in Discourse and Dialogue*. Vol. 39 of
18 *Text, Speech and Language Technology*. Springer.
- 19 Pietquin, O., 2004. A Framework for Unsupervised Learning of Dialogue
20 Strategies. SIMILAR Collection. Presses Universitaires de Louvain.
- 21 Pulman, S., 1996. Conversational games, belief revision and bayesian networks.
22 In: *Proceedings of the 7th Computational Linguistics in the Netherlands*
23 *meeting*.
- 24 Roy, N., Pineau, J., Thrun, S., 2000. Spoken Dialogue Management Using
25 Probabilistic Reasoning. In: *Proceedings of the Association for Computa-*
26 *tional Linguistics (ACL)*.
- 27 Schatzmann, J., 2008. Statistical user modeling for dialogue systems. Ph.D.
28 thesis, University of Cambridge.
- 29 Schatzmann, J., Thomson, B., Weillhammer, K., Ye, H., Young, S., 2007.
30 Agenda-based user simulation for bootstrapping a POMDP dialogue sys-
31 tem. In: *Proceedings of Human Language Technologies / North American*
32 *Chapter of the Association for Computational Linguistics (HLT/NAACL)*.
- 33 Scheffler, K., 2002. Automatic design of spoken dialogue systems. Ph.D. thesis,
34 University of Cambridge.
- 35 Shani, G., Poupart, P., Brafman, R., Shimony, S., 2008. Efficient add op-
36 erations for point-based algorithms. In: *The International Conference on*
37 *Automated Planning and Scheduling (ICAPS)*.
- 38 Sutton, R., Barto, A., 1998. *Reinforcement Learning: An Introduction*. Adap-
39 *tive Computation and Machine Learning*. MIT Press, Cambridge, Mass.
- 40 Sutton, R., McAllester, D., Singh, S., Mansour, Y., 2000. Policy gradient meth-
41 ods for reinforcement learning with function approximation. In: *Advances*
42 *in Neural Information Processing Systems 12*. MIT Press, pp. 1057–1063.
- 43 Thomson, B., Schatzmann, J., Weillhammer, K., Ye, H., Young, S., 2007.
44 Training a real-world POMDP-based dialog system. In: *Proceedings of the*
45 *HLT/NAACL workshop on “Bridging the Gap: Academic and Industrial*
46 *”*

- Research in Dialog Technologies". Rochester, NY.
- 1 Thomson, B., Yu, K., Gasic, M., Keizer, S., Mairesse, F., Schatzmann, J.,
2 Young, S., 2008. Evaluating semantic-level confidence scores with multiple
3 hypotheses. In: Proceedings of Interspeech. Brisbane, Australia.
- 4 Walker, M. A., 2000. An application of reinforcement learning to dialogue
5 strategy selection in a spoken dialogue system for email. *Journal of Artificial
6 Intelligence Research* 12, 387–416.
- 7
8 Williams, J., Young, S., 2006a. Scaling POMDPs for dialog management with
9 composite summary point-based value iteration (cspbvi). In: Proceedings of
10 the AAI Workshop on Statistical and Empirical Approaches for Spoken
11 Dialogue Systems. Boston.
- 12
13 Williams, J. D., 2007a. Applying POMDPs to dialog systems in the trou-
14 bleshooting domain. In: Proceedings of the HLT/NAACL workshop on
15 "Bridging the Gap: Academic and Industrial Research in Dialog Technol-
16 ogy". Rochester, NY, USA.
- 17
18 Williams, J. D., 2007b. Using particle filters to track dialogue state. In: IEEE
19 Workshop on Automatic Speech Recognition and Understanding (ASRU).
20 Kyoto, Japan.
- 21
22 Williams, J. D., Poupart, P., Young, S., 2005. Factored partially observable
23 Markov decision processes for dialogue management. In: Proceedings of the
24 IJCAI Workshop on Knowledge and Reasoning in Practical Dialog Systems.
25 Edinburgh.
- 26
27 Williams, J. D., Young, S., November 2005. Scaling up POMDPs for dialog
28 management: The "Summary POMDP" method. In: IEEE workshop on Au-
29 tomatic Speech Recognition and Understanding (ASRU). Cancun, Mexico.
- 30
31 Williams, J. D., Young, S., 2006b. Partially Observable Markov Decision Pro-
32 cesses for Spoken Dialog Systems. *Computer Speech and Language* 21 (2),
33 231–422.
- 34
35 Williams, J. D., Young, S., 2007. Scaling POMDPs for spoken dialog man-
36 agement. *IEEE Transactions on Audio, Speech, and Language Processing*
37 15(7), 2116–2129.
- 38
39 Yedidia, J. S., Freeman, W. T., Weiss, Y., 2001. Generalized belief propaga-
40 tion. In: *Advances in Neural Information Processing Systems* 13. MIT Press,
41 pp. 689–695.
- 42
43 Young, S., Schatzmann, J., Weilhammer, K., Ye, H., 2007. The Hidden Infor-
44 mation State Approach to Dialog Management. In: Proceedings of the Inter-
45 national Conference on Acoustics, Speech, and Signal Processing (ICASSP).
46 Honolulu, Hawaii.
- 47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65