

Recent Developments in Statistical Dialogue Systems

Steve Young



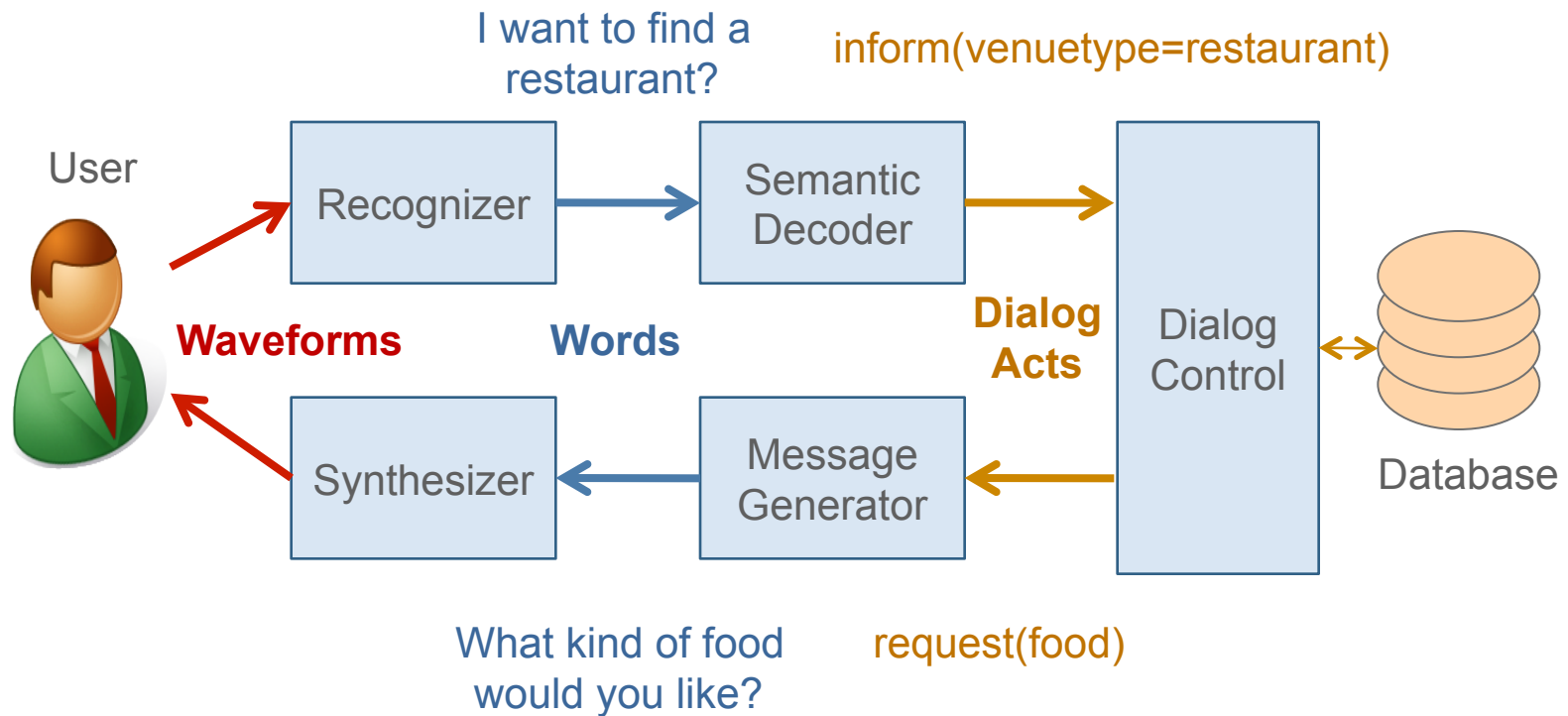
*Machine Intelligence Laboratory
Information Engineering Division
Cambridge University Engineering Department
Cambridge, UK*



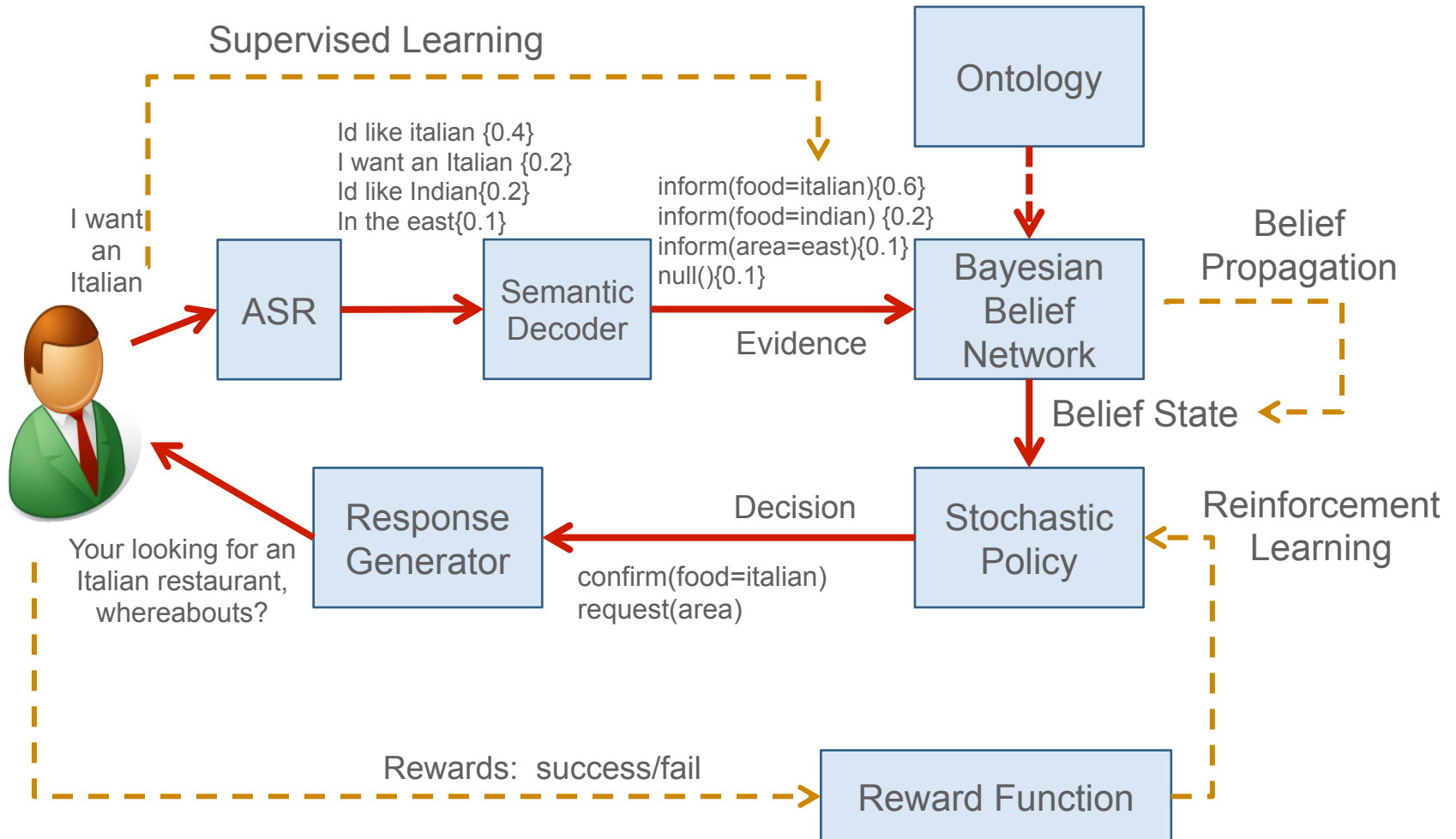


- Review of Basic Ideas and Current Limitations
- Semantic Decoding
- Fast Learning
- Parameter Optimisation and Structure Learning

Spoken Dialog Systems (SDS)

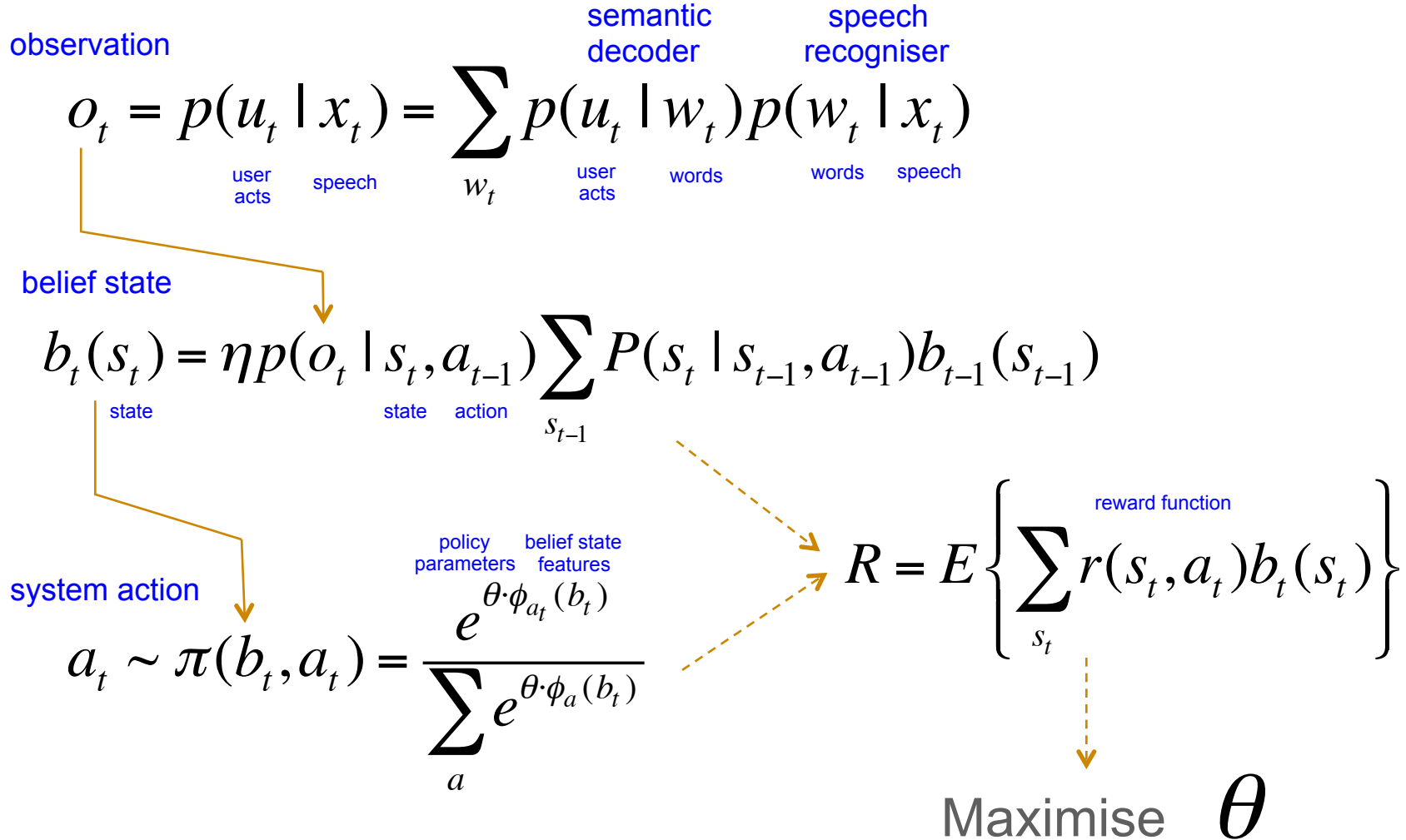


A Statistical Spoken Dialogue System



Partially Observable Markov Decision Process (POMDP)

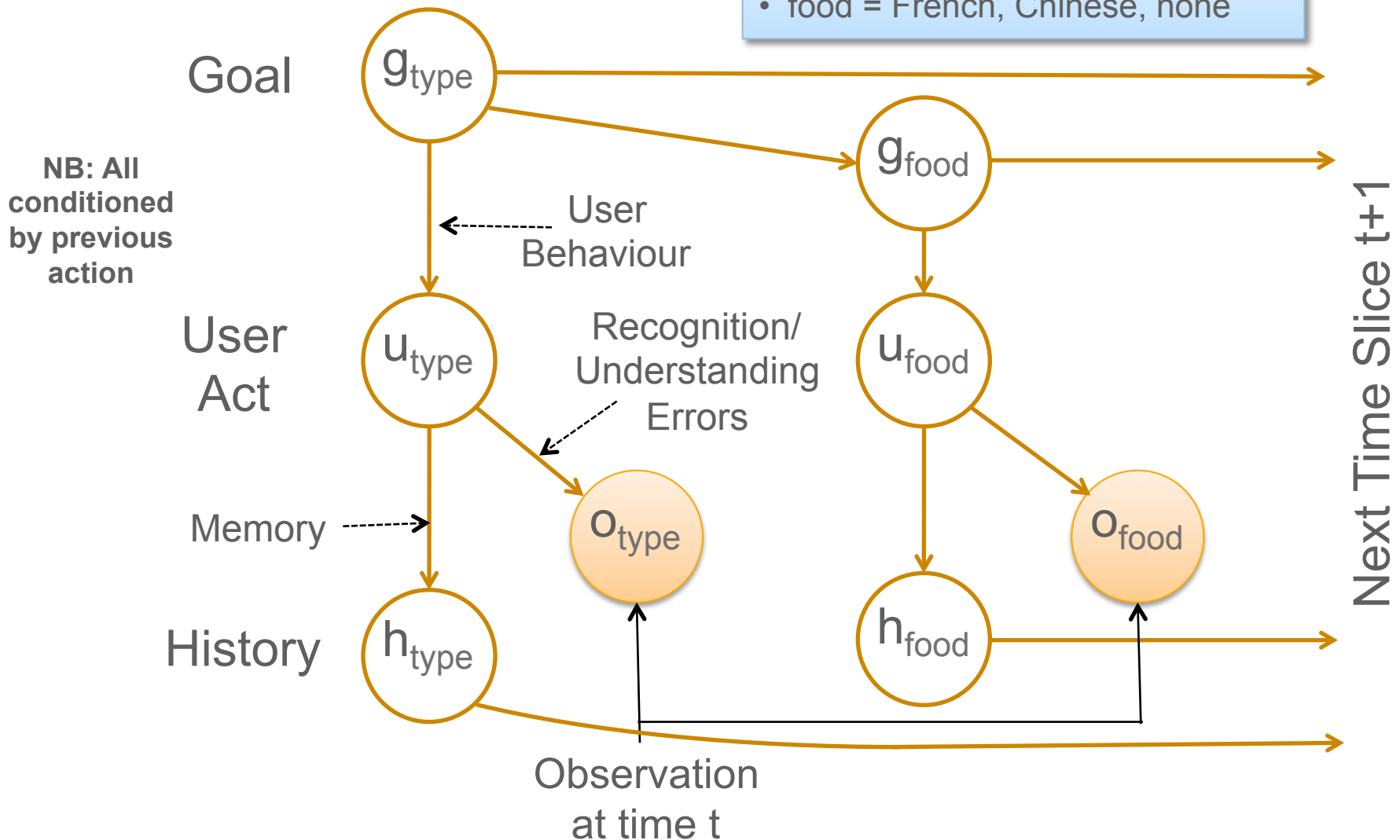
The POMDP SDS Framework



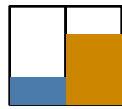
Dialogue State

Tourist Information Domain

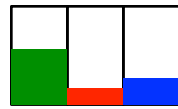
- type = bar, restaurant
- food = French, Chinese, none



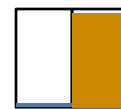
Belief Monitoring (Tracking)



B R



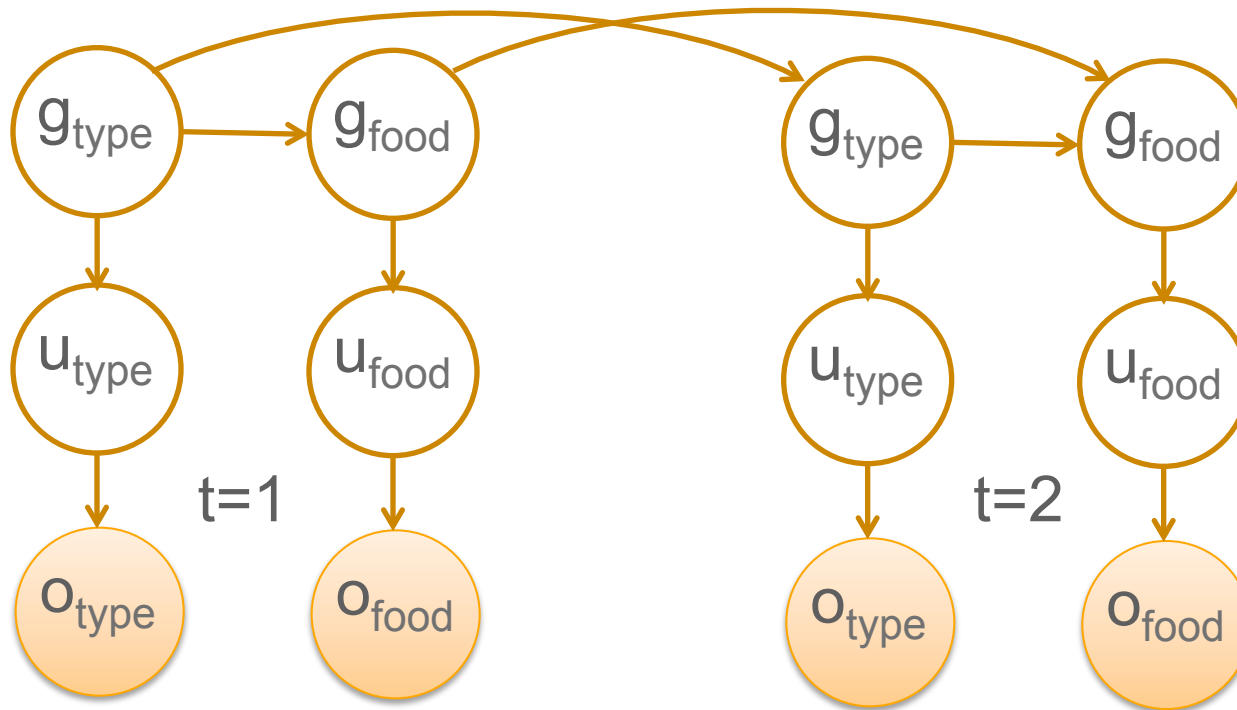
F C -



B R



F C -

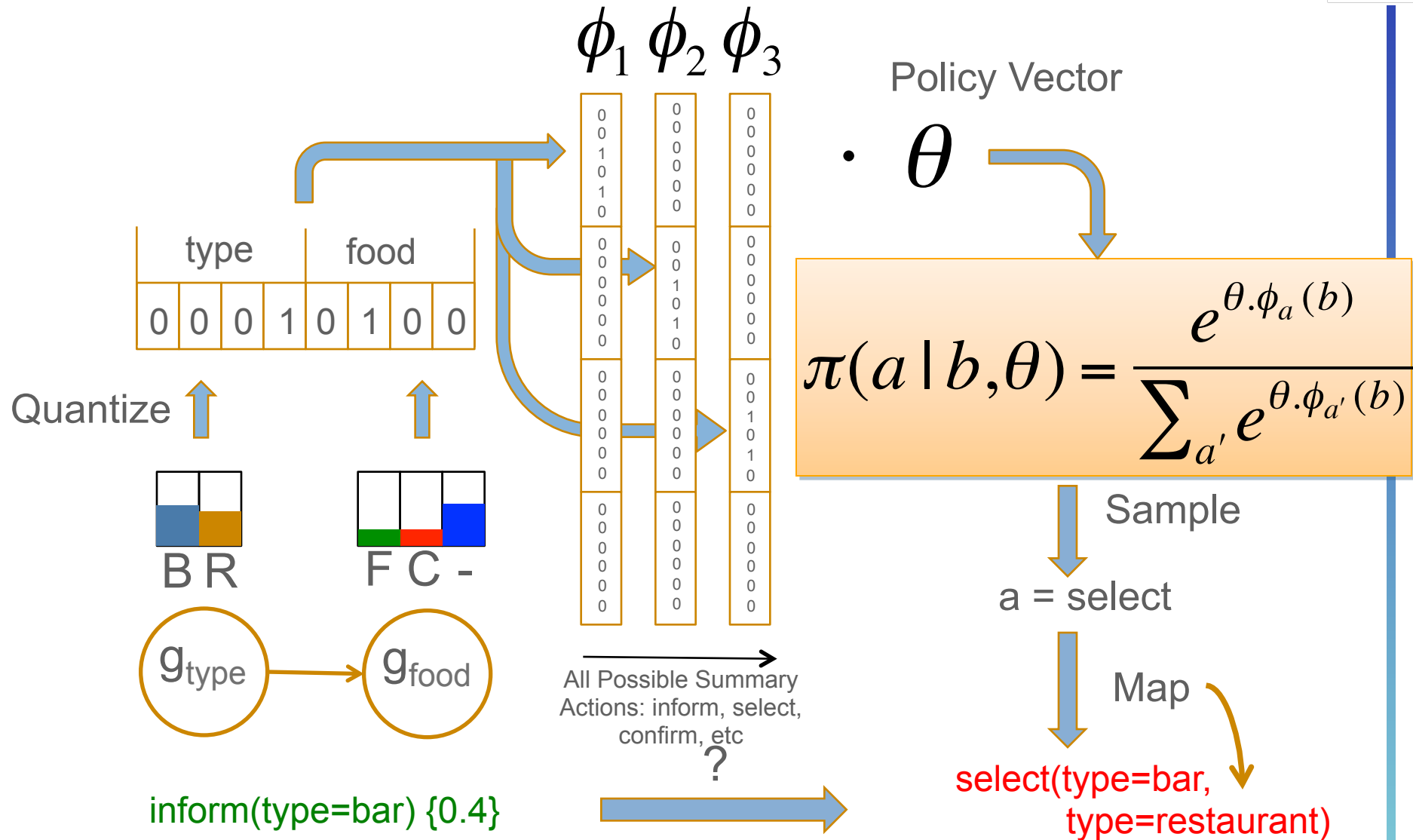


inform(type=bar,
food=french) {0.6}
inform(type=restaurant,
food=french) {0.3}

confirm(type=restaurant,
food=french)

affirm() {0.9}

Choosing the next action – the Policy

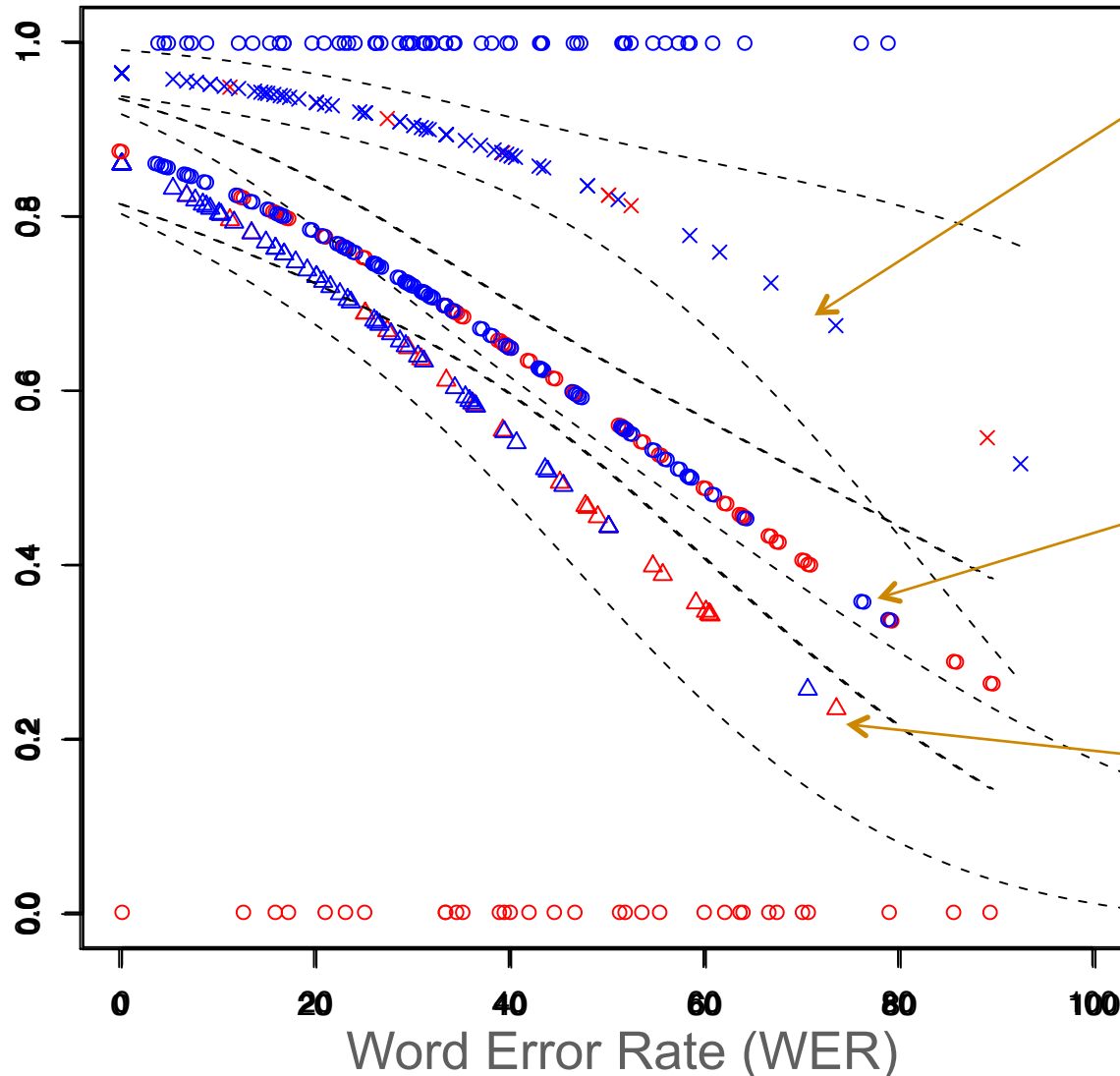


Let's Go 2010 Control Test Results



All Baseline Systems

Success
Success
Rate



Cambridge
89% Success
33% WER

Baseline
65% Success
42% WER

Another
75% Success
34% WER


B. Thomson
"Bayesian
Update of State
for the Let's Go
Spoken
Dialogue
Challenge."
SLT 2010.


Demo of Cambridge Restaurant Information




Call the system by pressing the call button to the right.



 Dialogue

 Map

 Details

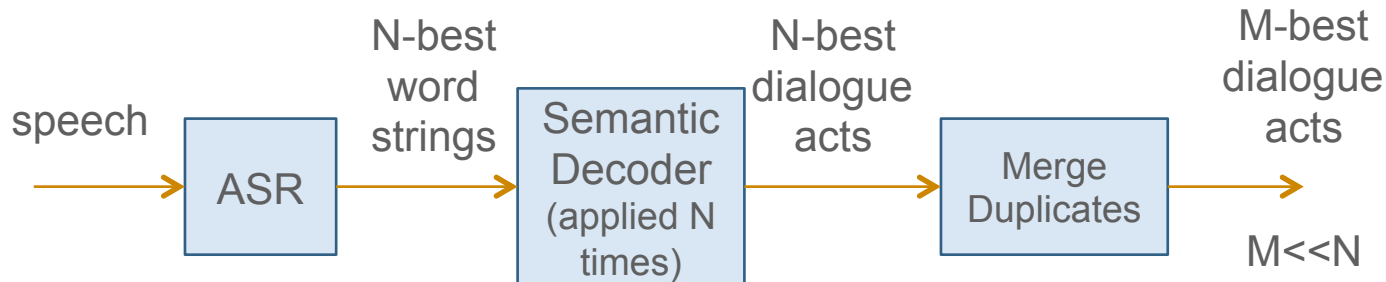


- Poor coverage of N-best of semantic hypotheses
- Hand-crafting of summary belief space
- Slow policy optimisation and reliance on user simulation
- Dependence on hand-crafted dialogue model parameters
- Dependence on static ontology/database

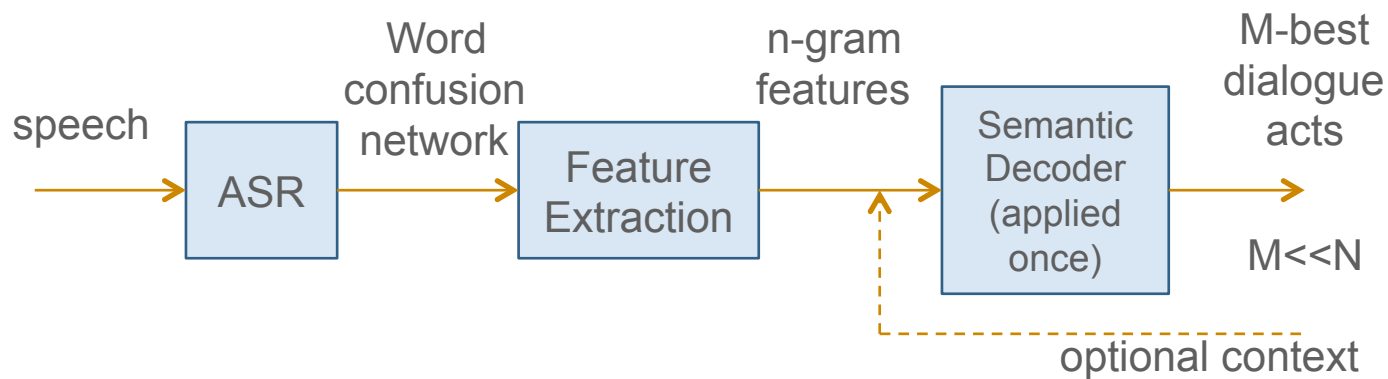
N-best Semantic Decoding



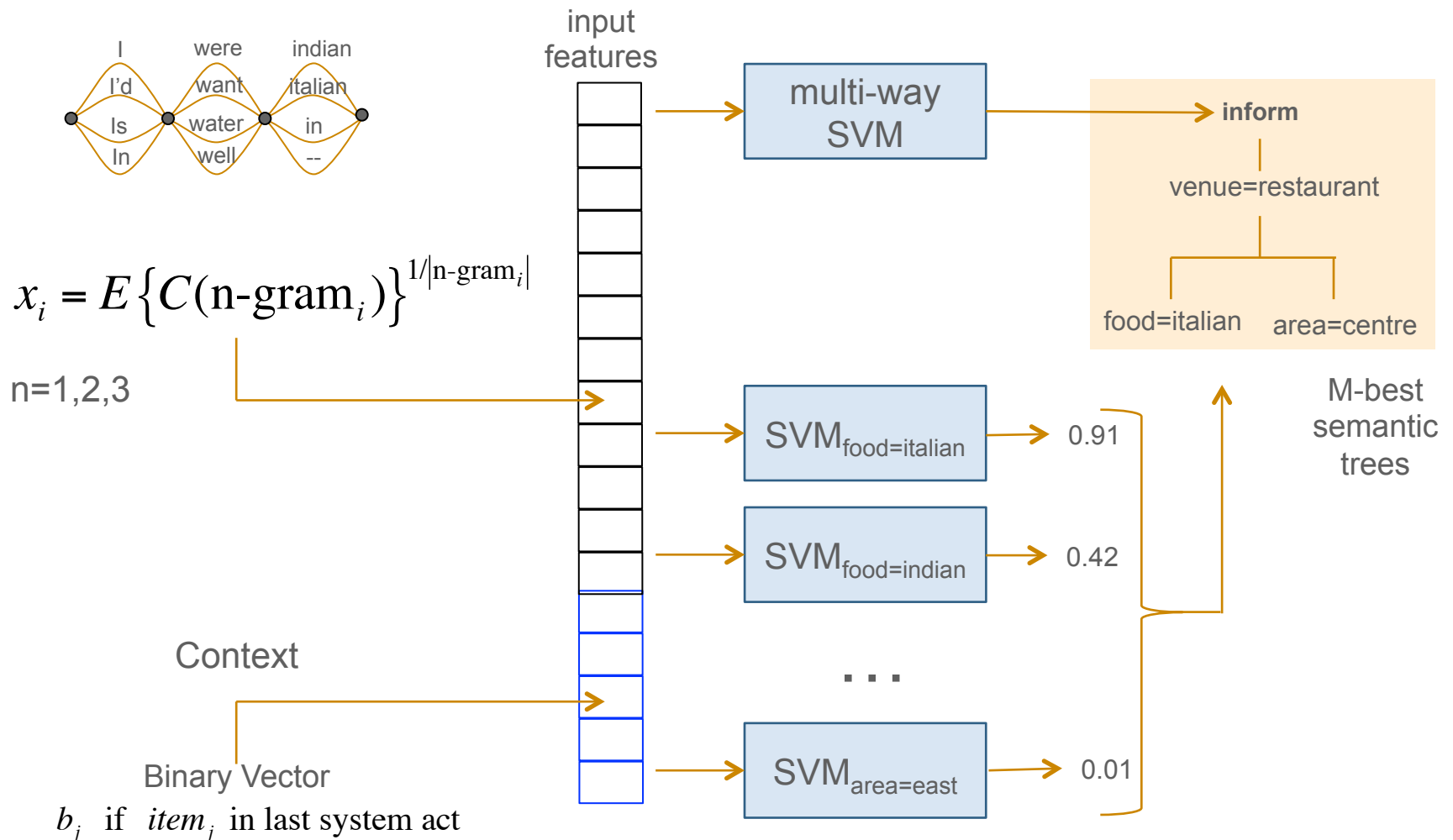
Conventional N-best decoding



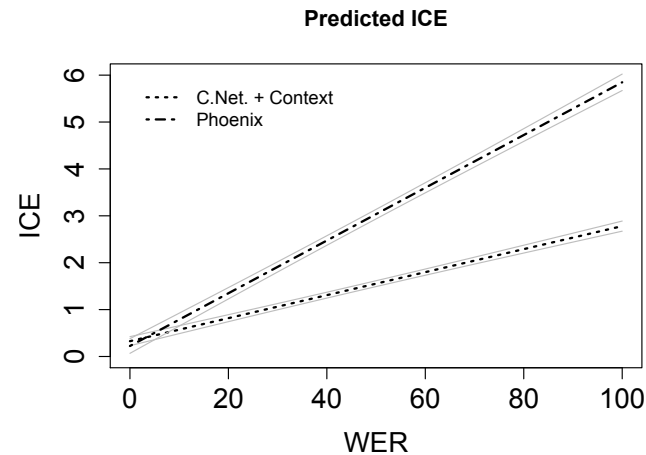
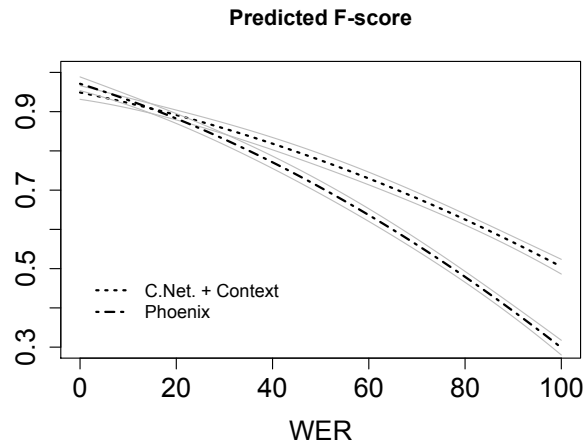
Confusion Network decoding



Confusion Network Decoder (Mairesse/Henderson)



Confusion Network Decoder Evaluation



Comparison of item retrieval on corpus of 4.8k utterances

N-best hand-crafted Phoenix decoder vs Confusion network decoder (trained on 10k utterances)

Live Dialogue System

	N-best Phoenix	Confusion Net
F-score	0.80	0.82
ICE	2.02	1.26
Average Reward	10.6	11.15

“Discriminative Spoken Language Understanding using Word Confusion Networks”, Henderson et al, IEEE SLT 2012



Policy parameters chosen to maximize expected reward

$$J(\theta) = E \left[\frac{1}{T} \sum_t r(s_t, a_t) \mid \pi_\theta \right]$$

Natural gradient ascent works well

$$\tilde{\nabla} J(\theta) = F_\theta^{-1} \nabla J(\theta)$$

Fisher
Information
Matrix



Gradient is estimated by sampling dialogues and in practice Fisher Information Matrix does not need to be explicitly computed. This is the Natural Actor Critic Algorithm.

However,

- A) slow (~100k dialogues) and
- B) requires summary space approximations

Q-functions and the SARSA algorithm



Traditional reinforcement learning is commonly based on finding the optimal Q function:

$$Q^*(b, a) = \max_{\pi} \left[E_{\pi} \left\{ \sum_{\tau=t+1}^T r(b_{\tau}, a_{\tau}) \right\} \right]$$

The optimal deterministic policy is then simply

$$\pi^*(b) = \operatorname{argmax}_a [Q^*(b, a)]$$

Q^* can be found sequentially using the SARSA algorithm

```
b = b0; choose action a e-greedily from π(b)
For each dialogue turn
    Take action a, observe reward r and next state b'
    choose action a' e-greedily from π(b')
    Q(b, a) = Q(b, a) + λ[Q(b', a') - (Q(b, a) - r)]
    b = b'; a = a'
end
```

Eventually, $Q \rightarrow Q^*$

Gaussian Process based Learning – Milica Gasic



For POMDPs, the belief space is continuous and direct representations of Q are intractable. However, Q can be approximated as a zero mean Gaussian process by designing a *kernel* to represent the correlations between points in belief \times action space. Thus:

$$Q(b, a) \sim GP(0, k((b, a), (b, a)))$$

Given a sequence of state-action pairs

$$B_t = [(b_0, a_0), (b_1, a_1), \dots, (b_t, a_t)]' \quad \text{and rewards} \quad r_t = [r_0, r_1, \dots, r_t]'$$

there is a closed form solution for the posterior:

$$Q(b, a) | B_t, r_t \sim N(\bar{Q}(b, a), \text{cov}((b, a), (b, a)))$$

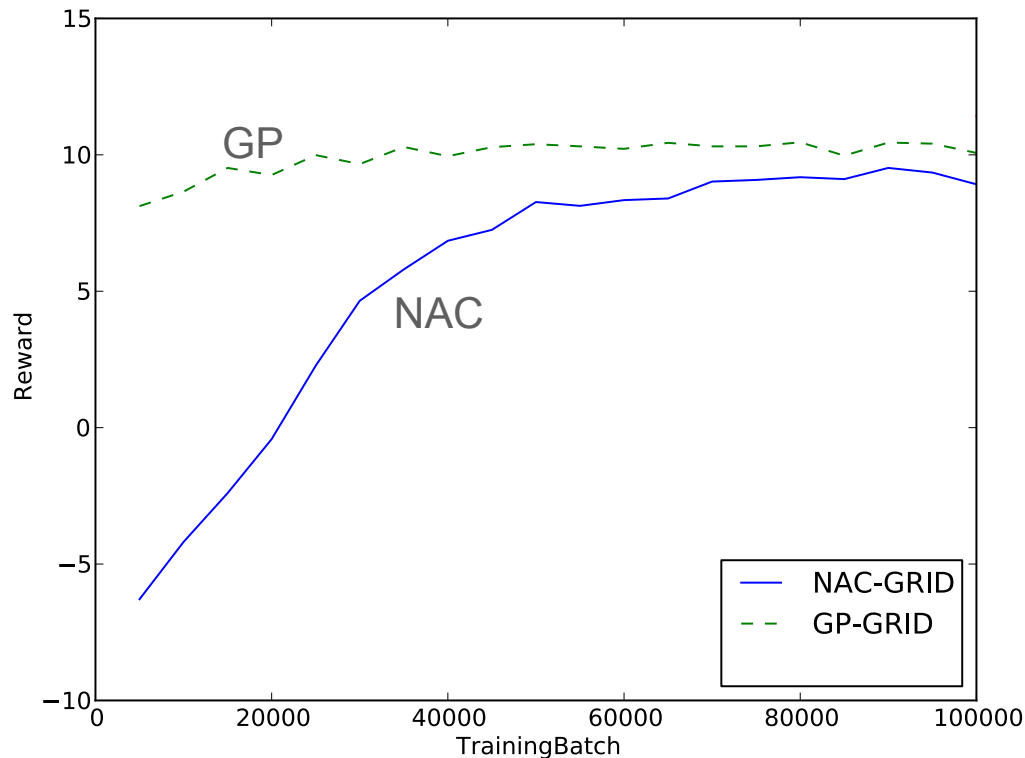
This suggests a SARSA-like sequential optimisation:

```
b = b0; choose action a e-greedily from  $\bar{Q}(b, a)$ 
For each dialogue turn
    Take action a, observe reward r and next state b'
    choose action a' e-greedily from  $\bar{Q}(b', a')$ 
    Update the posterior covariance estimate
    b = b'; a = a'
end
```

Benefits of GP-SARSA



- sequential estimation of distribution of Q (not Q itself)
- each new data point can impact on whole distribution via the covariance function
→ very efficient use of training data
- much faster learning than gradient methods such a Natural Actor Critic (NAC)



TownInfo
System

ϵ -greedy exploration

$$a = \begin{cases} \underset{a}{\operatorname{argmax}} [\bar{Q}(b, a)] & \text{with prob } 1 - \epsilon \\ \text{random action} & \text{with prob } \epsilon \end{cases}$$

Benefits of GP-SARSA



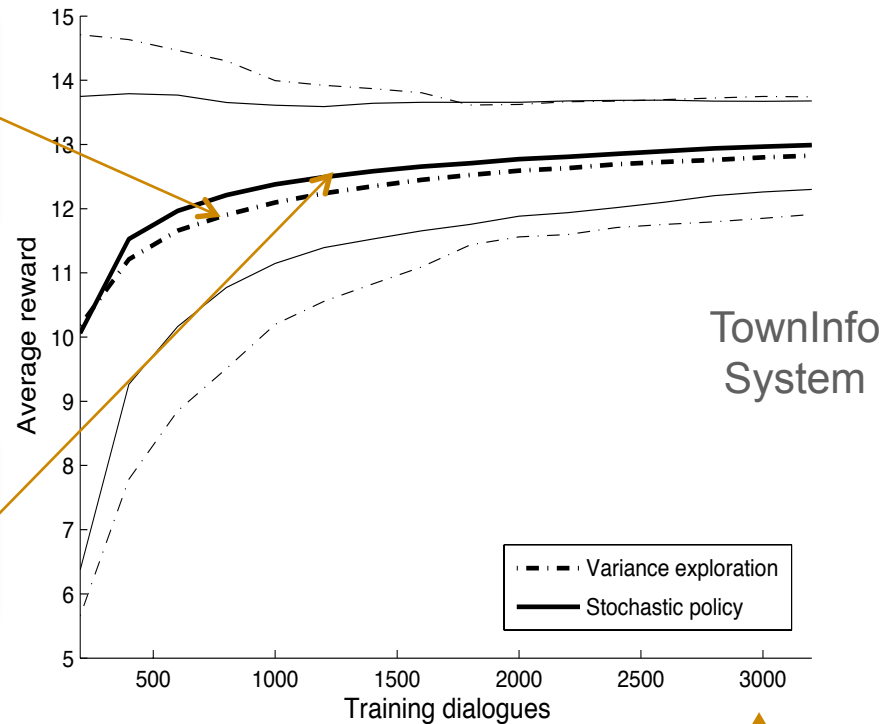
- variance of Q is known at each stage → more intelligent exploration:

Variance exploration

$$a = \begin{cases} \operatorname{argmax}_a [\bar{Q}(b, a)] & \text{with prob } 1 - \varepsilon \\ \operatorname{argmax}_a [\operatorname{cov}((b, a), (b, a))] & \text{with prob } \varepsilon \end{cases}$$

Stochastic policy

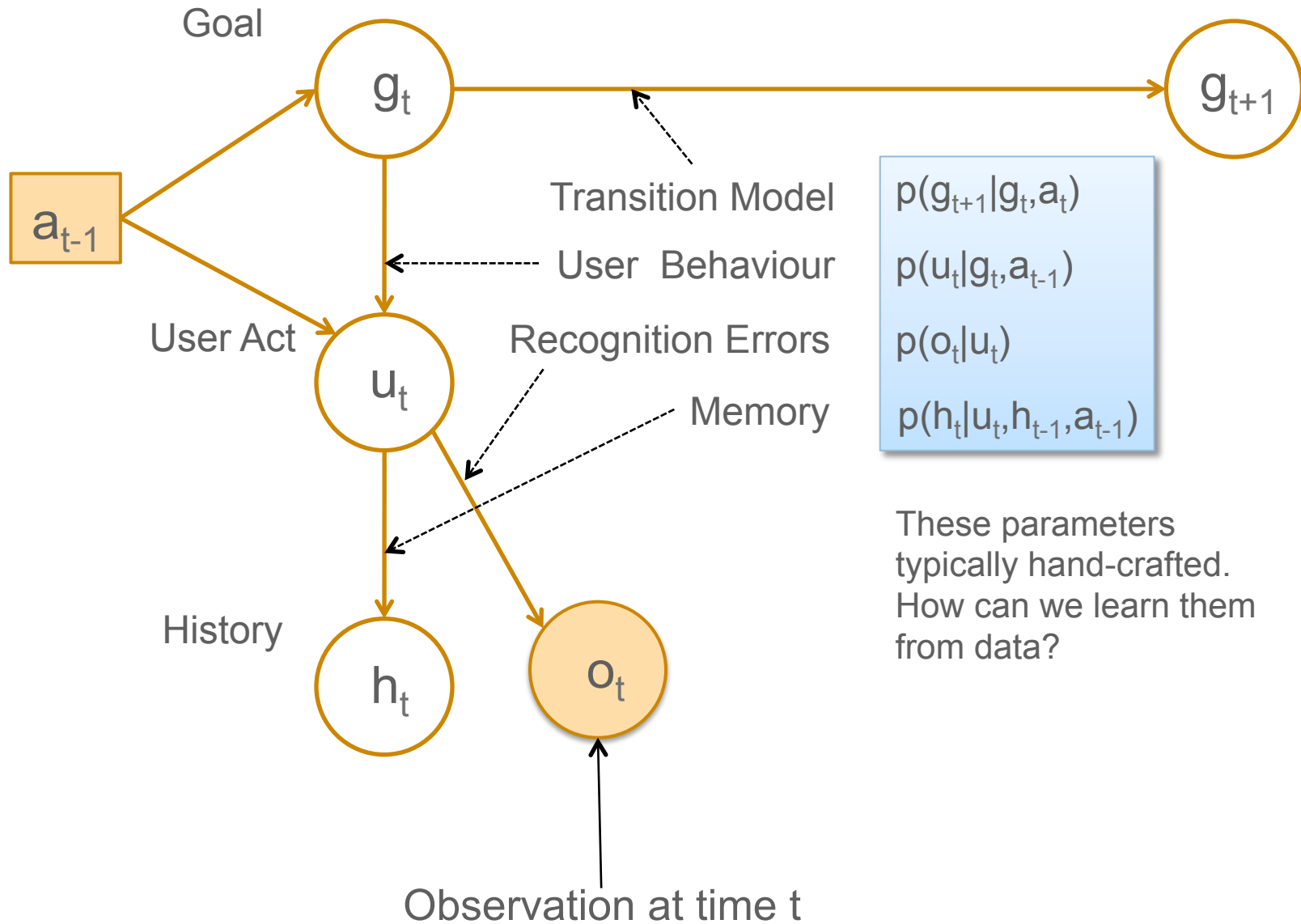
$$Q^i(b, a_i) \sim N(\bar{Q}(b, a_i), \operatorname{cov}((b, a_i), (b, a_i)))$$
$$a = \operatorname{argmax}_{a_i} [Q^i(b, a_i)]$$



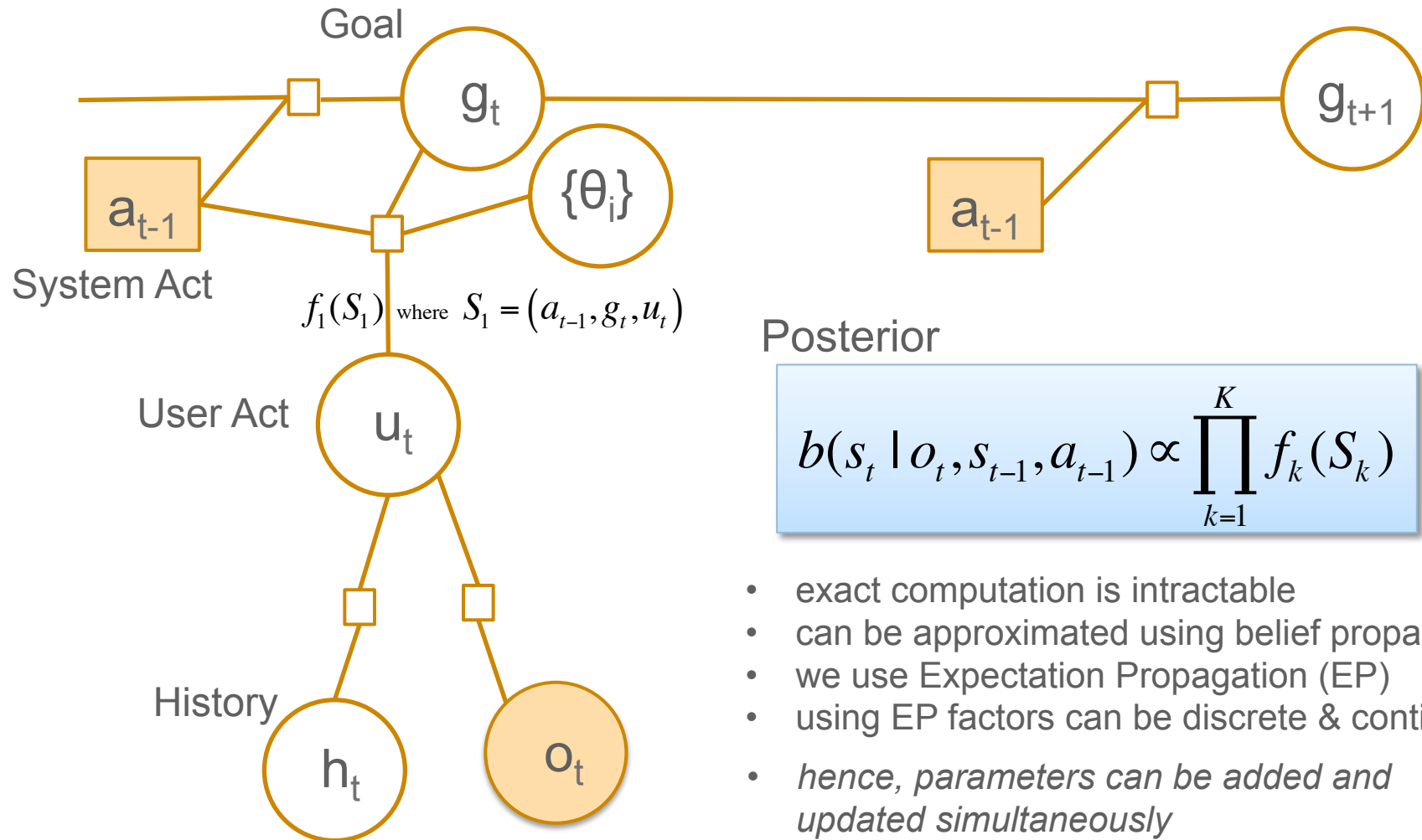
Well trained within 3k dialogues
And summary space mapping no longer needed

“On-line policy optimisation of SDS via live interaction with human subjects”, Gasic et al, ASRU 2011
“Gaussian processes for policy optimisation of large scale POMDP-based SDS”, Gasic et al, SLT 2012.

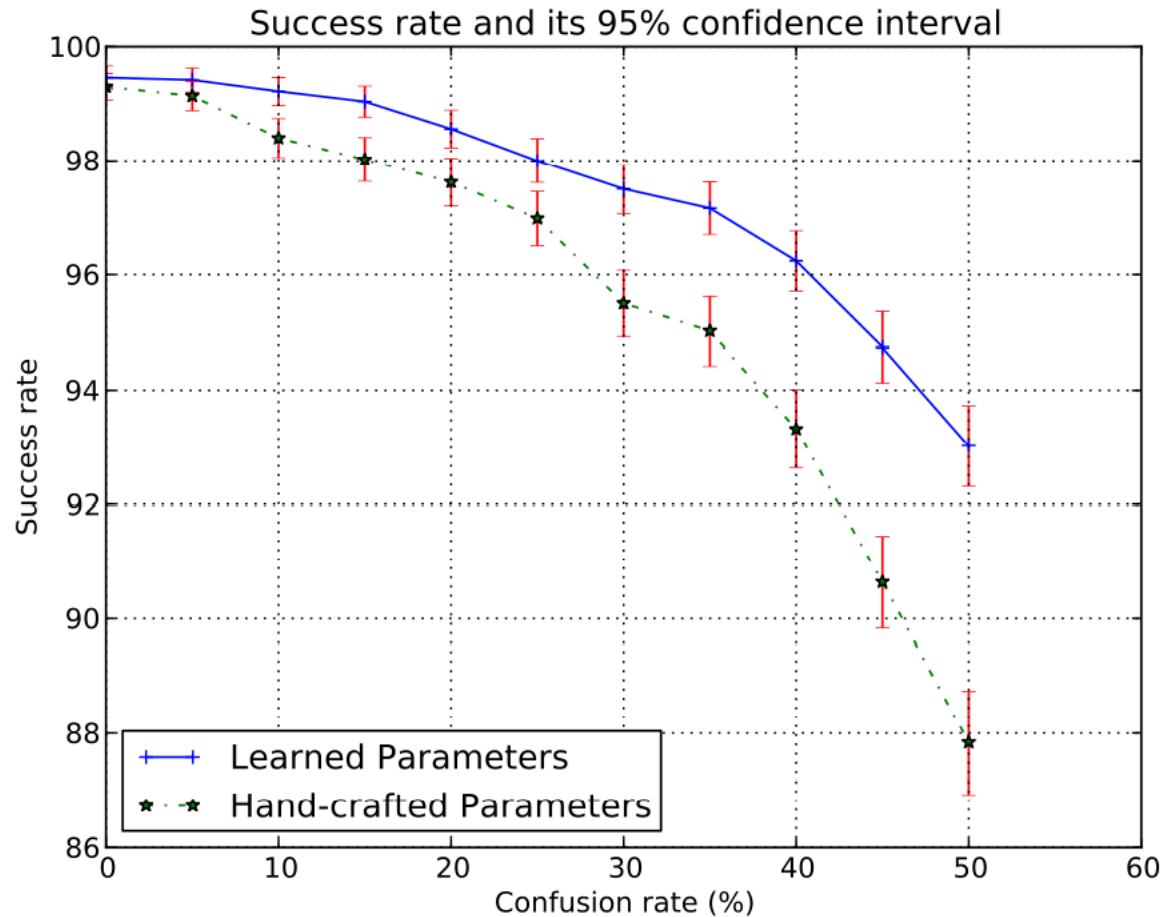
Parameter Estimation – Blaise Thomson



Factor Graphs and Expectation Propagation



Effect of Parameter Learning on TownInfo System

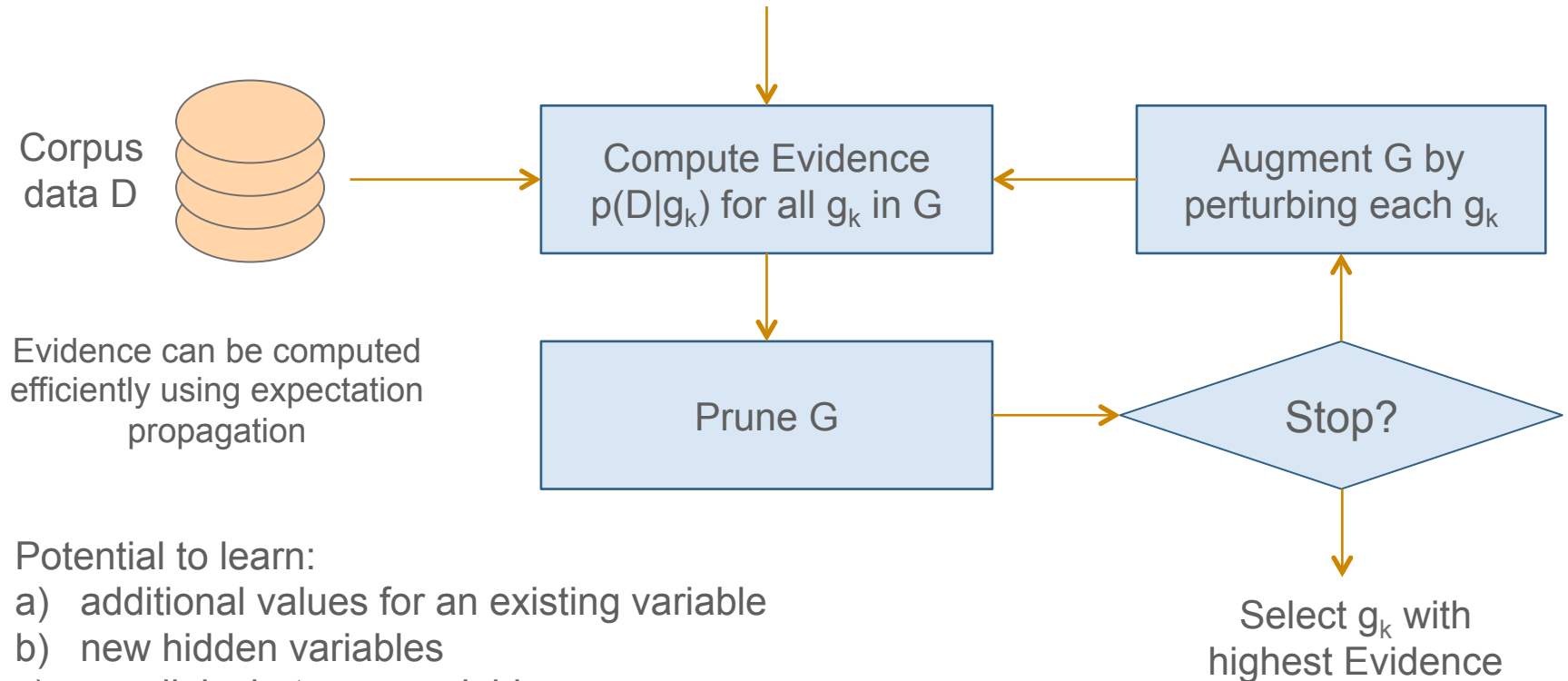


Structure Learning



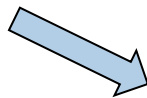
The ability to learn parameters can be extended to learn structure.

Let $G = \{g_k\}$ be a set of Bayesian Networks (or Factor Graphs):



Potential to learn:

- a) additional values for an existing variable
- b) new hidden variables
- c) new links between variables



PARLANCE



- Statistical Dialogue Systems based on POMDPs are viable, offer increased robustness to noise and require no hand-crafting
- Good progress is being made on increasing accuracy and speeding up learning
- Learning directly on human users rather than depending on user simulators is now possible
- Current systems are built from static ontologies for closed domains
- Next steps will include building more flexible systems capable of dynamically adapting to new information content.

Acknowledgement - this work is the result of a team effort: Catherin Breslin, Milica Gasic, Matt Henderson, Dongho Kim, Martin Szummer, Blaise Thomson, Pirros Tsiakoulis and former members of the CUED Dialogue Systems Group