

Just-In-Time Learning for Fast and Flexible Inference

S. M. Ali Eslami, Daniel Tarlow, Pushmeet Kohli and John Winn

{alie,dtarlow,pkohli,jwinn}@microsoft.com

Microsoft
Research

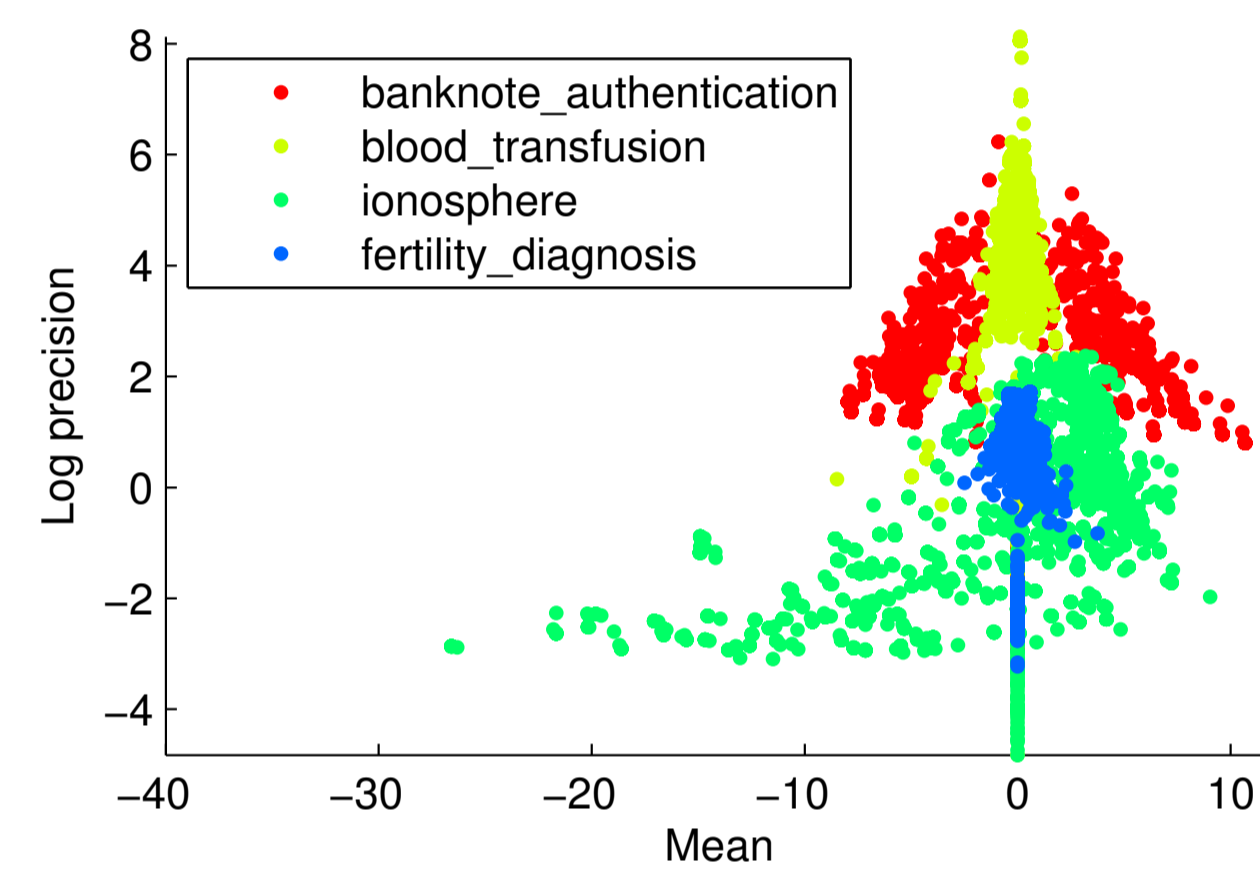
1. Goals

To combine the **flexibility** of sampling-based inference with the **speed** of optimized message-passing.

To allow practitioners to define any probabilistic model, press a button, and get accurate inference results within a matter of seconds.

2. Key observation

General algorithms appear to be solving problems that are harder than they need to be: in most **specific** inference problems, we only need to perform a small subset of all possible computations.



3. The JIT learning approach

- Initially use highly general algorithms for inference, e.g. by computing messages in a message-passing algorithm via sampling.
- Gradually learn to increase the speed of these computations by regressing from input to output messages at run-time.

4. JIT message computation

Learn a mapping from variable-to-factor messages $\{m_{k \rightarrow \psi}\}$ to a factor-to-variable message $m_{\psi \rightarrow i}$:

$$\bar{m}_{\psi \rightarrow i}(x_i) \equiv f(\{m_{k \rightarrow \psi}\}_{k \in \text{ne}(\psi)} | \theta). \quad (1)$$

Along with each prediction, produce a measure \bar{u} of the regressor's uncertainty:

$$\bar{u}_{\psi \rightarrow i} \equiv u(\{m_{k \rightarrow \psi}\}_{k \in \text{ne}(\psi)} | \theta). \quad (2)$$

Use uncertainties at run-time to choose between regressor estimates and slow 'oracle' computations:

$$m_{\psi \rightarrow i}(x_i) = \begin{cases} \bar{m}_{\psi \rightarrow i}(x_i) & \bar{u}_{\psi \rightarrow i} < u^{\max} \\ m_{\psi \rightarrow i}^{\text{oracle}}(x_i) & \text{otherwise} \end{cases} \quad (3)$$

where u^{\max} is the maximum tolerated uncertainty for a prediction. Set u^{\max} such that no held out prediction has an error above a user-specified maximum tolerated value D^{\max} .

Oracle messages can be computed using generally-applicable Monte Carlo (Barthelme *et al.*, 2011; Heess *et al.*, 2013), or via hand-implemented operators (e.g. those in Infer.NET, Minka *et al.*, 2012).

5. Random forests for JIT learning

Aim. Learn a mapping from incoming messages $\{m_{k \rightarrow \psi}\}_{k \in \text{ne}(\psi)}$ to the outgoing message $m_{\psi \rightarrow i}$.

Requirements. The regressor must: 1) train and predict efficiently, 2) model arbitrarily complex mappings, 3) adapt dynamically, and 4) produce uncertainty estimates.

Representation. Let m_{out} be the outgoing message and \mathbf{m}_{in} the set of incoming messages.

- Represent the outgoing message m_{out} by a vector of real valued numbers \mathbf{r}_{out} .
- Each set of incoming messages \mathbf{m}_{in} is represented in two ways:
 - Regression parameterization denoted by \mathbf{r}_{in} (e.g. concatenation of message parameters),
 - Tree parameterization denoted by \mathbf{t}_{in} (e.g. moments and ψ evaluated at the mode of \mathbf{m}_{in}).

Split and prediction model. \mathbf{t}_{in} is used to traverse message sets down to leaves, and \mathbf{r}_{in} is used by a polynomial regressor to predict \mathbf{r}_{out} :

$$\mathbf{r}_{\text{out}} = \mathbf{W} \cdot \phi^n(\mathbf{r}_{\text{in}}) + \epsilon. \quad (4)$$

Training objective function. At each node j , depending on the subset of the incoming training set \mathcal{S}_j we learn the function that 'best' splits \mathcal{S}_j into the training sets corresponding to each child, \mathcal{S}_j^L and \mathcal{S}_j^R , i.e. $\tau_j = \arg\max_{\tau \in \mathcal{T}_j} I(\mathcal{S}_j, \tau)$. The objective function I is:

$$I(\mathcal{S}_j, \tau) = -E(\mathcal{S}_j^L, \mathbf{W}^L) - E(\mathcal{S}_j^R, \mathbf{W}^R), \quad (5)$$

where \mathbf{W}^L and \mathbf{W}^R are the parameters of the polynomial regression models corresponding to the left and right training sets \mathcal{S}_j^L and \mathcal{S}_j^R , and the 'fit residual' E is:

$$E(\mathcal{S}, \mathbf{W}) = \frac{1}{2} \sum_{\mathbf{m}_{\text{in}} \in \mathcal{S}} D_{\text{KL}}^{\text{mar}}(\bar{\mathbf{m}}_{\mathbf{m}_{\text{in}}}^{\mathbf{W}} \| m_{\mathbf{m}_{\text{in}}}^{\text{oracle}}) + D_{\text{KL}}^{\text{mar}}(m_{\mathbf{m}_{\text{in}}}^{\text{oracle}} \| \bar{\mathbf{m}}_{\mathbf{m}_{\text{in}}}^{\mathbf{W}}). \quad (6)$$

This objective function splits the training data at each node in a way that the relationship between the incoming and outgoing messages is well captured by the polynomial regression in each child.

Mean squared error of message parameters is sensitive to parameterization. Instead, calculate the marginals \bar{b}_i and b_i^{oracle} induced on the target variable, and compute their KL:

$$D_{\text{KL}}^{\text{mar}}(\bar{m}_{\psi \rightarrow i} \| m_{\psi \rightarrow i}^{\text{oracle}}) \equiv D_{\text{KL}}(\bar{b}_i \| b_i^{\text{oracle}}), \quad (7)$$

$D_{\text{KL}}^{\text{mar}}$ is the *marginal KL* and is used throughout the JIT framework, as it encourages the system to focus efforts on the quantity that is ultimately of interest: the accuracy of the posterior marginals.

Ensemble model. Compute the *moment average* \bar{m}_{out} of the distributions $\{\bar{m}_{\text{out}}^t\}$ by averaging the first few moments of each predicted distribution across trees, and solving for the distribution parameters which match the averaged moments.

The moment average can be interpreted as minimizing an objective function:

$$\bar{m}_{\text{out}} = \underset{m}{\text{argmin}} U(\{\bar{m}_{\text{out}}^t\}, m), \quad (8)$$

where

$$U(\{\bar{m}_{\text{out}}^t\}, m) = \sum_t D_{\text{KL}}(\bar{m}_{\text{out}}^t \| m). \quad (9)$$

The level of agreement between the predictions of the different trees can be used as a proxy of the forest's uncertainty about that prediction. If all the trees in the forest predict the same output distribution, it means that their knowledge about the function f is similar *despite* the randomness in their structures. We therefore set $\bar{u}_{\text{out}} \equiv U(\{\bar{m}_{\text{out}}^t\}, \bar{m}_{\text{out}})$.

6. Experimental results

Logistic factor

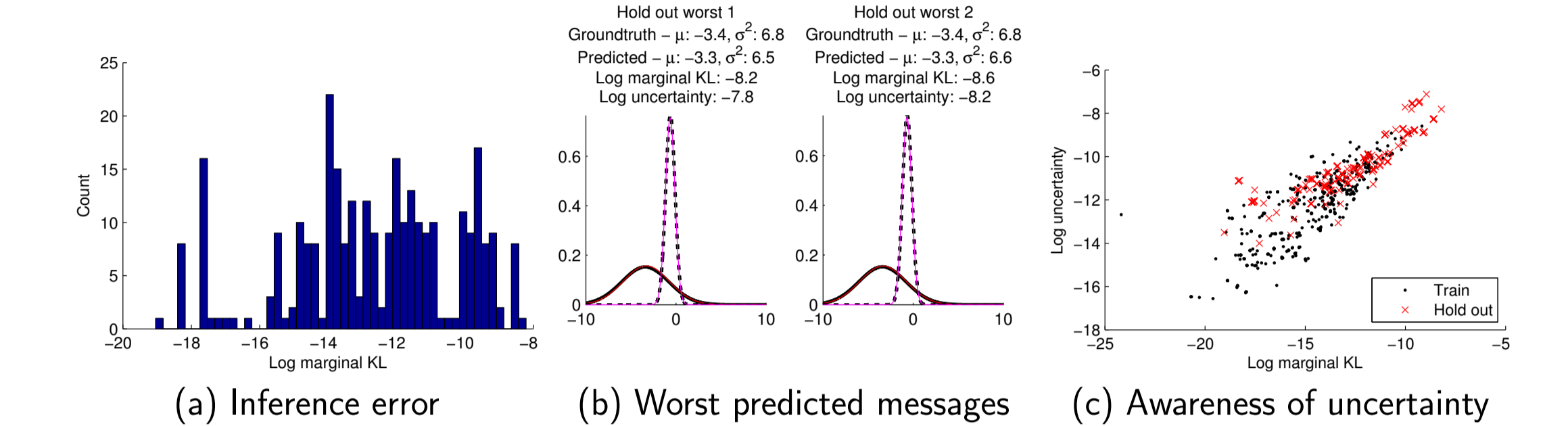


Figure 1: **Uncertainty aware regression.** Learning a logistic factor for use in logistic regression. (a) Histogram of marginal KLs of outgoing messages (Gaussian), which are typically very small. (b) The forest's most inaccurate predictions (black: m^{oracle} , red: \bar{m} , dashed black: b^{oracle} , purple: \bar{b}). (c) The regressor's uncertainty increases with marginal KL, i.e. it does not make confident but inaccurate predictions.

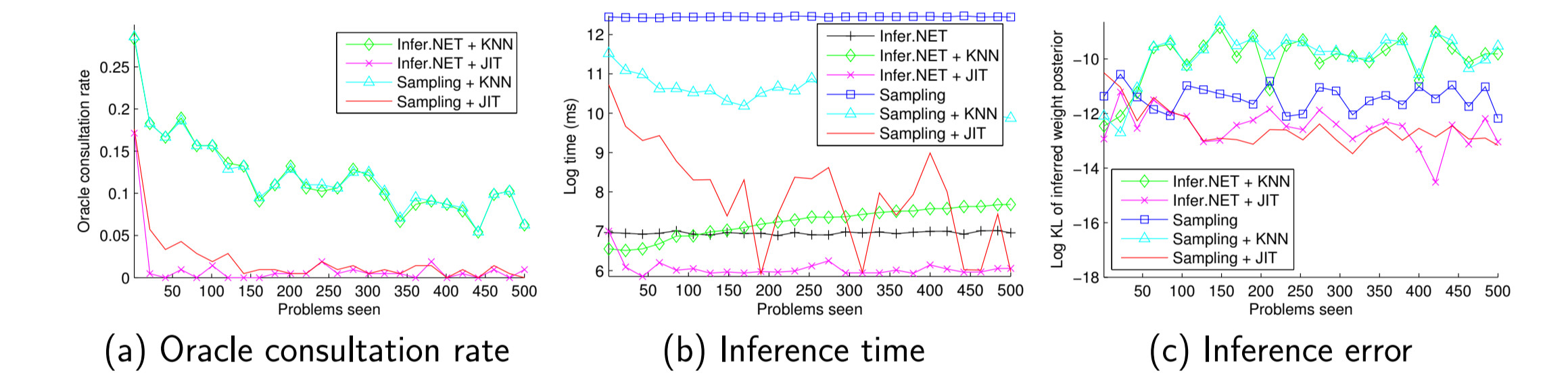


Figure 2: (a) The factor consults the oracle for only a fraction of messages, (b) leading to significant savings in time, (c) whilst maintaining (or decreasing) error.

Compound gamma factor

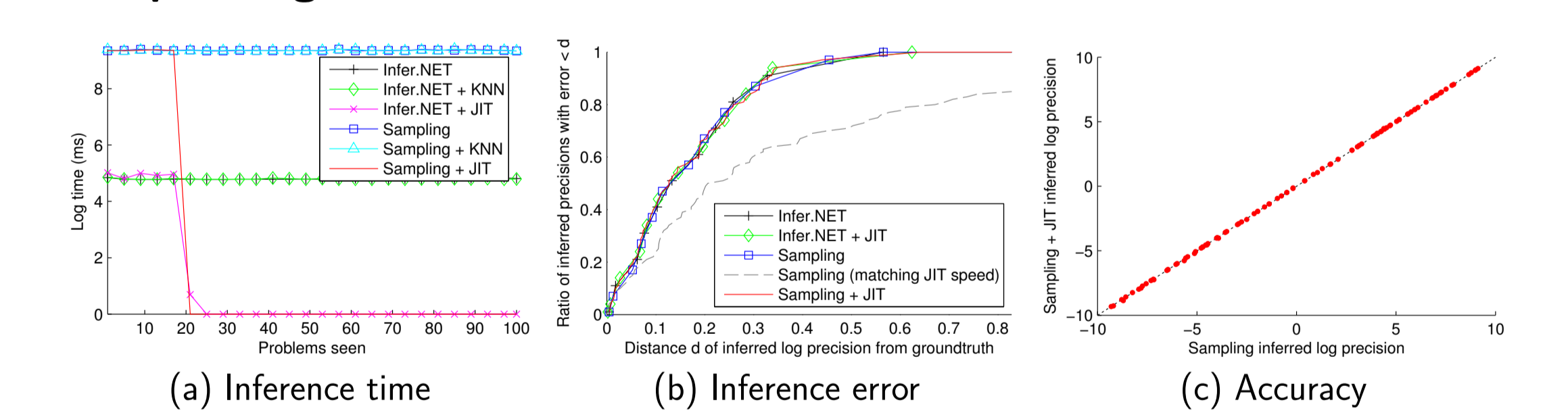
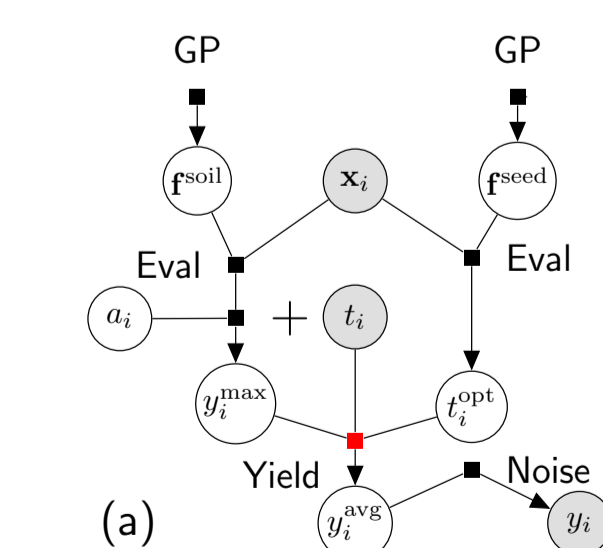


Figure 3: (a) JIT reduces inference time for sampling from ~ 11 seconds to ~ 1 ms. (b) JIT's posteriors agree highly with Infer.NET. Using fewer samples to match JIT speed leads to degradation of accuracy. (c) Speed comes at limited loss of accuracy.

Yield factor



(b)	IS		JIT fresh		JIT continued	
	Time	FR	Speedup	FR	Speedup	Speedup
11	451s	54%	195%	—	—	—
12	449s	54%	192%	60%	288%	—
13	451s	54%	191%	64%	318%	—

Figure 4: (a) The yield factor relates temperatures and yields recorded at farms to the optimal temperatures of their planted grain. JIT learning enables us to incorporate arbitrary factors with ease, whilst maintaining inference speed. (b) FR is fraction of regressions with no oracle consultation.