

---

# Continuous-Time Belief Propagation

---

Tal El-Hay<sup>1</sup>  
Ido Cohn<sup>1</sup>  
Nir Friedman<sup>1</sup>  
Raz Kupferman<sup>2</sup>

TALE@CS.HUJI.AC.IL  
IDO.COHN@CS.HUJI.AC.IL  
NIR@CS.HUJI.AC.IL  
RAZ@MATH.HUJI.AC.IL

<sup>1</sup>School of Computer Science and Engineering, <sup>2</sup>Institute of Mathematics, Hebrew University, Jerusalem 91904, Israel

## Abstract

Many temporal processes can be naturally modeled as a stochastic system that evolves continuously over time. The representation language of *continuous-time Bayesian networks* allows to succinctly describe multi-component continuous-time stochastic processes. A crucial element in applications of such models is (approximate) inference. Here we introduce a variational approximation scheme, which is a natural extension of Belief Propagation for continuous-time processes. In this scheme, we view messages as inhomogeneous Markov processes over individual components. This leads to a relatively simple procedure that allows to easily incorporate adaptive ordinary differential equation (ODE) solvers to perform individual steps. We provide the theoretical foundations for the approximation, and show how it performs on a range of networks. Our results demonstrate that our method is quite accurate on singly connected networks, and provides close approximations in more complex ones.

## 1. Introduction

The dynamics of many real-life processes are naturally modeled in terms of continuous-time stochastic processes, allowing for a wide range of time scales within the same process. Examples include biological sequence evolution (Felsenstein, 2004), computer systems (Xu & Shelton, 2008; Simma et al., 2008), and social networks (Fan & Shelton, 2009).

While the mathematical foundations of continuous-time stochastic processes are well understood (Chung, 1960),

the study of efficient computer representations, inference, and learning of complex continuous-time processes is still in early stages. *Continuous-time Bayesian networks* (CTBNs) (Nodelman et al., 2002) provide a sparse representation of complex multi-component processes by describing how the dynamics of an individual component depends on the state of its neighbors. A major challenge is translating the structure of a CTBN to computational gains in inference problems—answering queries about the process from partial observations.

As exact inference in a CTBN is exponential in the number of components, we have to resort to approximations. Broadly speaking, these fall into two main categories. The first category is stochastic approximations (Fan & Shelton, 2008; El-Hay et al., 2008), which sample trajectories of the process. While these can be asymptotically exact, they can be computationally expensive and incur computational penalties when sampling rapidly evolving processes. The second category of approximations is variational methods. Nodelman et al. (2005) and Saria et al. (2007) developed an approach based on *expectation propagation* (Minka, 2001; Heskes & Zoeter, 2002), where the posterior distribution over a process is approximated by *piecewise homogeneous* factored processes. This involves an elaborate message passing scheme between the approximations for different components, and an adaptive procedure for determining how to segment each time interval. More recently, Cohn et al. (2009) introduced a mean-field approximation (Jordan et al., 1998), which uses factored *inhomogeneous* processes (Opper & Sanguinetti, 2007). This allowed them to build on the rich literature of adaptive ODE solvers. While the mean-field approximation provides a lower-bound on the likelihood, it suffers from the expected drawbacks when approximating highly coupled processes.

Here we introduce a variational approximation that combines insights from both previous approaches for variational inference in CTBNs. Our approximation is a natural extension of the successful Bethe approximation (Yedidia et al., 2005) to CTBNs. Alternatively, it can be viewed

---

Appearing in *Proceedings of the 27<sup>th</sup> International Conference on Machine Learning*, Haifa, Israel, 2010. Copyright 2010 by the author(s)/owner(s).

applying the approach of Nodelman et al where the segment length diminishes to zero. Our approximation finds a collection of inhomogeneous processes over subsets of components, which are constrained to be locally consistent over single components. We show that this approximation is often accurate on tree-networks, and provides good approximations for more complex networks. Importantly, the approximation scheme is simple and allows to easily exploit the large suites of computational tools offered in the field of ODEs.

## 2. Continuous-Time Bayesian Networks

Consider a  $d$ -component Markov process  $\mathbf{X}^{(t)} = (X_1^{(t)}, X_2^{(t)}, \dots, X_d^{(t)})$  with state space  $S = S_1 \times S_2 \times \dots \times S_d$ . A notational convention: vectors are denoted by bold-face symbols, e.g.,  $\mathbf{X}$ , and matrices are denoted by black-board style characters, e.g.,  $\mathbb{Q}$ . The states in  $S$  are denoted by vectors of indexes,  $\mathbf{x} = (x_1, \dots, x_d)$ . We use indexes  $1 \leq i, j \leq d$  for enumerating components and  $\mathbf{X}^{(t)}$  and  $X_i^{(t)}$  to denote the random variable describing the state of the process and its  $i$ 'th component at time  $t$ .

The dynamics of a *time-homogeneous continuous-time Markov process* are fully determined by the *Markov transition function*,

$$p_{\mathbf{x}, \mathbf{y}}(t) = \Pr(\mathbf{X}^{(t+s)} = \mathbf{y} | \mathbf{X}^{(s)} = \mathbf{x}),$$

where time-homogeneity implies that the right-hand side does not depend on  $s$ . These dynamics are captured by a matrix  $\mathbb{Q}$ —the *rate matrix*, with non-negative off-diagonal entries  $q_{\mathbf{x}, \mathbf{y}}$  and diagonal entries  $q_{\mathbf{x}, \mathbf{x}} = -\sum_{\mathbf{y} \neq \mathbf{x}} q_{\mathbf{x}, \mathbf{y}}$ . The rate matrix is related to the transition function by

$$\left. \frac{d}{dt} p_{\mathbf{x}, \mathbf{y}}(t) \right|_{t=0} = q_{\mathbf{x}, \mathbf{y}}.$$

The probability of being in state  $\mathbf{x}$  at time  $t$  satisfies the *master equation* (Chung, 1960)

$$\frac{d}{dt} \Pr(\mathbf{X}^{(t)} = \mathbf{x}) = \sum_{\mathbf{y}} q_{\mathbf{y}, \mathbf{x}} \Pr(\mathbf{X}^{(t)} = \mathbf{y}).$$

A *continuous-time Bayesian network* is a structured multi-component continuous-time Markov process. It is defined by assigning each component  $i$  a set of components  $\mathbf{Pa}_i \subseteq \{1, \dots, d\} \setminus \{i\}$ , which are its parents in the network (Nodelman et al., 2002). With each component  $i$  we then associate a family of rate matrices  $\mathbb{Q}_{\mathbf{u}_i}^{i|\mathbf{Pa}_i}$ , with entries  $q_{x_i, y_i}^{i|\mathbf{Pa}_i}$ , that describe the rates of change of the  $i$ 'th component given the state  $\mathbf{u}_i$  of the parents  $\mathbf{Pa}_i$ . The dynamics of  $\mathbf{X}^{(t)}$  are defined by a rate matrix  $\mathbb{Q}$  with entries  $q_{\mathbf{x}, \mathbf{y}}$

that combines the conditional rate matrices as follows:

$$q_{\mathbf{x}, \mathbf{y}} = \begin{cases} q_{x_i, y_i}^{i|\mathbf{Pa}_i} & \delta_{\mathbf{x}, \mathbf{y}} = \{i\} \\ \sum_i q_{x_i, x_i}^{i|\mathbf{Pa}_i} & \mathbf{x} = \mathbf{y} \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

where  $\delta_{\mathbf{x}, \mathbf{y}} = \{i | x_i \neq y_i\}$ . This definition implies that changes occur one component at a time.

Given a continuous-time Bayesian network, we would like to evaluate the likelihood of evidence, to compute the probability of various events given the evidence (e.g., that the state of the system at time  $t$  is  $\mathbf{x}$ ), and to compute conditional expectations (e.g., the expected amount of time  $X_i$  was in state  $x_i$ ). Direct computations of these quantities involve matrix exponentials of the rate matrix  $\mathbb{Q}$ , whose size is exponential in the number of components, making this approach infeasible beyond a modest number of components. We therefore have to resort to approximations.

## 3. A Variational Principle

Variational inference methods pose the inference task in terms of an optimization problem. The objective is to maximize a functional which lower-bounds the log probability of the evidence by introducing an auxiliary set of *variational parameters* (Wainwright & Jordan, 2008). Recently, Cohn et al. (2009) introduced a variational formulation of inference in continuous-time Markov processes. We start by reviewing the relevant results of Cohn et al.

For convenience we restrict our treatment to an interval  $[0, T]$  with boundary evidence  $\mathbf{X}^{(0)} = \mathbf{e}_0$  and  $\mathbf{X}^{(T)} = \mathbf{e}_T$ . The posterior distribution of a homogeneous Markov process given evidence  $\mathbf{e} = \{\mathbf{e}_0, \mathbf{e}_T\}$  on the two boundaries is a *non-homogeneous Markov process*. Such a process can be represented using a *time varying rate matrix*  $\mathbb{Q}(t)$  that describe the instantaneous transition rates. However, such a representation is unwieldy, since as  $t$  approaches  $T$  the transition rates from  $\mathbf{x} \neq \mathbf{e}_T$  to  $\mathbf{e}_T$  approach infinity.

To circumvent the problem of unbounded values near the boundaries, Cohn et al introduce *marginal density sets* which represent the posterior process in terms of univariate and joint pairwise distributions. More formally, if  $\Pr$  denotes the posterior distribution, its *marginal density set* is the following family of continuous functions:

$$\begin{aligned} \mu_{\mathbf{x}}(t) &= \Pr(\mathbf{X}^{(t)} = \mathbf{x}) \\ \gamma_{\mathbf{x}, \mathbf{y}}(t) &= \lim_{h \downarrow 0} \frac{\Pr(\mathbf{X}^{(t)} = \mathbf{x}, \mathbf{X}^{(t+h)} = \mathbf{y})}{h}, \quad \mathbf{x} \neq \mathbf{y}. \end{aligned} \quad (2)$$

In addition to providing a bounded representation to the posterior, this representation allows to easily compute ex-

pected sufficient statistics using numerical integration:

$$\mathbf{E}[T_{\mathbf{x}}(t)] = \int_0^t \mu_{\mathbf{x}}(s) ds, \quad \mathbf{E}[M_{\mathbf{x},\mathbf{y}}(t)] = \int_0^t \gamma_{\mathbf{x},\mathbf{y}}(s) ds,$$

where  $T_{\mathbf{x}}(t)$  is the residence time in state  $\mathbf{x}$  in the interval  $[0, t]$ , and  $M_{\mathbf{x},\mathbf{y}}(t)$  is the number of transitions from  $\mathbf{x}$  to  $\mathbf{y}$  in the same interval. Thus, this representation is analogous to sets of *mean parameters* that are employed in variational approximations over exponential families with a finite dimensional parametrization (Wainwright & Jordan, 2008; Koller & Friedman, 2009).

Families of functions  $\mu, \gamma$  that satisfy (2) for some Pr, must satisfy self-consistent relations imposed by the master equation.

**Definition 3.1:** (Cohn et al., 2009) A family  $\eta = \{\mu_{\mathbf{x}}(t), \gamma_{\mathbf{x},\mathbf{y}}(t) : 0 \leq t \leq T\}$  of continuous functions is a *Markov-consistent density set* if the following constraints are fulfilled:

$$\begin{aligned} \mu_{\mathbf{x}}(t) &\geq 0, \quad \sum_{\mathbf{x}} \mu_{\mathbf{x}}(0) = 1, \\ \gamma_{\mathbf{x},\mathbf{y}}(t) &\geq 0 \quad \forall \mathbf{y} \neq \mathbf{x}, \\ \frac{d}{dt} \mu_{\mathbf{x}}(t) &= \sum_{\mathbf{y} \neq \mathbf{x}} (\gamma_{\mathbf{y},\mathbf{x}}(t) - \gamma_{\mathbf{x},\mathbf{y}}(t)). \end{aligned}$$

and  $\gamma_{\mathbf{x},\mathbf{y}}(t) = 0$  whenever  $\mu_{\mathbf{x}}(t) = 0$ . For convenience, we define  $\gamma_{\mathbf{x},\mathbf{x}} = -\sum_{\mathbf{y} \neq \mathbf{x}} \gamma_{\mathbf{x},\mathbf{y}}$ . ■

The evidence at the boundaries impose additional constraints on potential posterior processes. Specifically, the representation  $\eta$  corresponding to the posterior distribution  $P_{\mathbb{Q}}(\cdot | e_0, e_T)$  is in the set  $\mathcal{M}_e$  that contains Markov-consistent density sets  $\{\mu_{\mathbf{x}}(t), \gamma_{\mathbf{x},\mathbf{y}}(t)\}$ , that satisfy  $\mu_{\mathbf{x}}(0) = \mathbf{1}_{\mathbf{x}=e_0}$ ,  $\mu_{\mathbf{x}}(T) = \mathbf{1}_{\mathbf{x}=e_T}$  and  $\gamma_{\mathbf{x},\mathbf{y}}(T) = 0$  for all  $\mathbf{y} \neq e_T$ . In addition, since these sets are posteriors of a CTBN, they also change one component at a time, which implies that  $\gamma_{\mathbf{x},\mathbf{y}}(t) = 0$  if  $|\delta_{\mathbf{x},\mathbf{y}}| > 1$ .

Using this representation, the variational formulation of Cohn et al is reminiscent of similar formulations for discrete probabilistic models (Jordan et al., 1998).

**Theorem 3.2:** (Cohn et al., 2009) Let  $\mathbb{Q}$  be a rate matrix and  $e = (e_0, e_T)$  be states of  $\mathbf{X}$ . Then

$$\ln P_{\mathbb{Q}}(e_T | e_0) = \max_{\eta \in \mathcal{M}_e} \mathcal{F}(\eta; \mathbb{Q}),$$

where

$$\mathcal{F}(\eta; \mathbb{Q}) = \mathcal{E}(\eta; \mathbb{Q}) + \mathcal{H}(\eta),$$

is the *free energy functional* which is a sum of an *average energy functional*

$$\mathcal{E}(\eta; \mathbb{Q}) = \int_0^T \sum_{\mathbf{x}} \left[ \mu_{\mathbf{x}}(t) q_{\mathbf{x},\mathbf{x}} + \sum_{\mathbf{y} \neq \mathbf{x}} \gamma_{\mathbf{x},\mathbf{y}}(t) \ln q_{\mathbf{x},\mathbf{y}} \right] dt,$$

and an *entropy functional*

$$\mathcal{H}(\eta) = \int_0^T \sum_{\mathbf{x}} \sum_{\mathbf{y} \neq \mathbf{x}} \gamma_{\mathbf{x},\mathbf{y}}(t) [1 + \ln \mu_{\mathbf{x}}(t) - \ln \gamma_{\mathbf{x},\mathbf{y}}(t)] dt.$$

To illustrate this principle, we can examine how to derive an exact inference procedure. We can find the optimum of  $\mathcal{F}(\eta; \mathbb{Q})$  by introducing Lagrange multipliers that enforce the consistency constraint, and then find the stationary point of the corresponding Lagrangian. Since we are dealing with a continuous-time formula, we need to use the Euler-Lagrange method (Gelfand & Fomin, 1963). As Cohn et al. (2009) show, the maximum satisfies a system of differential equations:

$$\begin{aligned} \frac{d}{dt} \rho_{\mathbf{x}} &= - \sum_{\mathbf{y}} q_{\mathbf{x},\mathbf{y}} \rho_{\mathbf{y}} & \rho_{\mathbf{x}}(T) &= \mathbf{1}_{\mathbf{x}=e_T} \\ \frac{d}{dt} \mu_{\mathbf{x}} &= \sum_{\mathbf{y} \neq \mathbf{x}} (\gamma_{\mathbf{y},\mathbf{x}} - \gamma_{\mathbf{x},\mathbf{y}}), & \mu_{\mathbf{x}}(0) &= \mathbf{1}_{\mathbf{x}=e_0} \\ \gamma_{\mathbf{x},\mathbf{y}} &= \mu_{\mathbf{x}} q_{\mathbf{x},\mathbf{y}} \frac{\rho_{\mathbf{y}}}{\rho_{\mathbf{x}}}, & \mathbf{y} \neq \mathbf{x}, \rho_{\mathbf{x}} &\neq 0, \end{aligned} \quad (3)$$

where we omit the  $(t)$  argument for clarity. The auxiliary functions  $\rho_{\mathbf{x}}(t)$  are Lagrange multipliers.

These equations have a simple intuitive solution that involves backward integration of  $\rho_{\mathbf{x}}(t)$ , as we have a boundary condition at time  $T$  and  $\rho_{\mathbf{x}}(t)$  does not depend on  $\mu_{\mathbf{x}}(t)$ . This integration results in

$$\rho_{\mathbf{x}}(t) = \Pr(e_T | \mathbf{X}^{(t)} = \mathbf{x})$$

Once we solve for  $\rho_{\mathbf{x}}(t)$ , we can forward integrate  $\mu_{\mathbf{x}}(t)$  from the boundary conditions at 0 to get the solution for  $\mu_{\mathbf{x}}$  and  $\gamma_{\mathbf{x},\mathbf{y}}$ . This analysis suggests that this system of ODEs is similar to forward-backward propagation, except that unlike classical forward propagation, here the forward propagation takes into account the backward messages to directly compute the posterior. Note that applying this exact solution to a multi-component process results in an exponential (in  $d$ ) number of coupled differential equations.

## 4. Continuous-Time Expectation Propagation

Approximate variational inference procedures are derived by posing an optimization problem that is an approximate version of the original one. Different approximations differ in terms of whether they approximate the objectives, limit or relax the allowed set of solutions, or combine several such approaches. Here, we follow a strategy which is based on the approach of *expectation propagation*, in which the set of admissible solutions is extended to ones that are consistent only on the expectations of statistics of interest, and in addition, use an approximate functional.

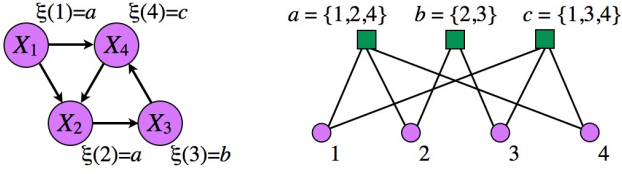


Figure 1. A CTBN and a corresponding factor graph.

#### 4.1. Approximate Optimization Problem

To represent potential solutions, we follow methods used in recent approximate inference procedures that use factor graph representations (Yedidia et al., 2005; Koller & Friedman, 2009). Specifically, we keep only marginal density sets over smaller clusters of components.

We start with definitions and notations. A *factor graph* is an undirected bipartite graph. One layer in the graph consists of *component nodes* that are labeled by component indexes. The second layer consists of *clusters nodes*  $\mathcal{A}$ , where each cluster  $\alpha \in \mathcal{A}$ , is a subset of  $\{1, \dots, d\}$ . The edges in the graph are between a component node  $i$  to a cluster node  $\alpha$  if and only if  $i \in \alpha$ . Thus, the neighbors of  $\alpha$  are  $N(\alpha) = \{i : i \in \alpha\}$  and the neighbors of  $i$  are  $N(i) = \{\alpha : i \in \alpha\}$ .

A factor graph is *family preserving*, with respect to a given CTBN, if there exists an assignment function  $\xi(i)$  that maps components to clusters, such that for every  $i$ , we have that  $\{i\} \cup \mathbf{Pa}_i \subseteq \xi(i)$ . We denote by  $A(\alpha)$  the set of components  $i$  for which  $\xi(i) = \alpha$ . From now on, we assume that we deal only with family preserving factor graphs.

**Example 4.1:** Figure 1 shows a simple CTBN and a corresponding factor graph. In this specific factor graph,  $\mathcal{A}(a) = \{1, 2\}$ ,  $\mathcal{A}(b) = \{3\}$  and  $\mathcal{A}(c) = \{4\}$ . ■

Given a factor graph, we use its structure to define an approximation for a distribution. Instead of describing the distribution over all the components, we use a family of density sets  $\tilde{\eta} = \{\eta^i : i = 1, \dots, d\} \cup \{\eta^\alpha : \alpha \in \mathcal{A}\}$ . A family of marginal density sets can be inconsistent. We do not require full consistency, but only consistency between neighboring nodes in the following sense.

**Definition 4.2:** A family of density sets  $\tilde{\eta}$  is said to be *locally consistent* if for all  $\alpha \in \mathcal{A}$  and all  $i \in N(\alpha)$  we have  $\mu^i = \mu^\alpha|_i$  where

$$(\mu^\alpha|_i)_{x_i} = \sum_{\mathbf{x}_{\alpha \setminus i}} \mu_{[\mathbf{x}_{\alpha \setminus i}, x_i]}^\alpha \quad (4)$$

and  $[\mathbf{x}_{\alpha \setminus i}, x_i]$  is the assignment to  $\mathbf{x}_\alpha$  composed from  $\mathbf{x}_{\alpha \setminus i}$

and  $x_i$ . Likewise,  $\gamma^i = \gamma^\alpha|_i$  where

$$(\gamma^\alpha|_i)_{x_i, y_i} = \sum_{\mathbf{x}_{\alpha \setminus i}} \gamma_{[\mathbf{x}_{\alpha \setminus i}, x_i], [\mathbf{x}_{\alpha \setminus i}, y_i]}^\alpha. \quad (5)$$

Let  $\tilde{\mathcal{M}}_e$  be the set of locally consistent densities that correspond to evidence  $e$  ■

The local consistency of  $\eta^\alpha$  and  $\eta^i$  does not imply that the distribution  $\text{Pr}_{\eta^i}(X_i)$  is equal to the marginal distribution  $\text{Pr}_{\eta^\alpha}(X_i)$ , as marginalization of a Markov process is not necessarily a Markov process. Rather,  $\text{Pr}_{\eta^i}$  is the projection of  $\text{Pr}_{\eta^\alpha}(X_i)$  to a Markov process with the matching expectations of  $\mathbf{E}[T_{x_i}(t)]$  and  $\mathbf{E}[M_{x_i, y_i}(t)]$  (Koller & Friedman, 2009).

Such locally consistent sets allow us to construct a tractable approximation to the variational optimization problem by introducing the *continuous-time Bethe functional*

$$\tilde{\mathcal{F}}(\tilde{\eta}; \mathbb{Q}) = \sum_i \mathcal{E}_i(\eta^{\alpha(i)}; \mathbb{Q}^{i|\mathbf{Pa}_i}) + \sum_\alpha \mathcal{H}(\eta^\alpha) - \sum_i c_i \mathcal{H}(\eta^i)$$

where

$$\mathcal{E}_i(\eta^\alpha; \mathbb{Q}^{i|\mathbf{Pa}_i}) = \int_0^T \sum_{\mathbf{x}_\alpha} \left[ \mu_{\mathbf{x}_\alpha}^\alpha(t) q_{x_i, x_i | \mathbf{u}_i}^{i|\mathbf{Pa}_i} + \sum_{\mathbf{y} \neq \mathbf{x}} \gamma_{\mathbf{y}, \mathbf{x}}^\alpha(t) \ln q_{x_i, y_i | \mathbf{u}_i}^{i|\mathbf{Pa}_i} \right] dt,$$

and  $c_i = N(i) - 1$  ensure that the total weight of sets containing component  $i$  sums up to 1. This functional is analogous to the well-known Bethe approximation for discrete models (Yedidia et al., 2005).

Combining the two approximations the approximate optimization problem becomes:

$$\max_{\tilde{\eta} \in \tilde{\mathcal{M}}_e} \tilde{\mathcal{F}}(\tilde{\eta}; \mathbb{Q}) \quad (6)$$

Once the optimal parameters are found, we can use the relevant marginal density set to answer queries.

#### 4.2. Stationary Point Characterization

To characterize the stationary points of the approximate optimization problem (6) we use again the Euler-Lagrange method, where we introduce Lagrange multiplier functions to enforce the cluster-wise constraints,  $\frac{d}{dt} \mu_{\mathbf{x}_\alpha}^\alpha = \sum_{\mathbf{y} \neq \mathbf{x}} (\gamma_{\mathbf{y}, \mathbf{x}}^\alpha - \gamma_{\mathbf{x}, \mathbf{y}}^\alpha)$  as well as the local consistency constraints defined in equations (4) and (5). Differentiating the Lagrangian, equating the derivatives to zero, and performing some algebra, which we omit for the lack of space, we obtain fixed-point equations that consist of the initial constraints and two classes of coupled equations.

The first class consists of equations similar to (3), which refer to the dynamics within each cluster. To simplify the presentation, we introduce some definitions.

**Definition 4.3:** Assume we are given a time-varying matrix function  $\mathbb{G}(t)$ , and boundary conditions  $\mathbf{x}_0$  and  $\mathbf{x}_T$ . Define the operator  $\eta = \mathcal{R}(\mathbb{G}, \mathbf{x}_0, \mathbf{x}_T)$  to return  $\eta = (\mu, \gamma)$ , the unique solution of the following ODEs

$$\begin{aligned} \frac{d}{dt} \rho_{\mathbf{x}} &= - \sum_{\mathbf{y}} g_{\mathbf{x}, \mathbf{y}} \rho_{\mathbf{y}}, & \rho_{\mathbf{x}}(T) &= \mathbf{1}_{\mathbf{x}=\mathbf{x}_T} \\ \frac{d}{dt} \mu_{\mathbf{x}} &= \sum_{\mathbf{y} \neq \mathbf{x}} (\gamma_{\mathbf{y}, \mathbf{x}} - \gamma_{\mathbf{x}, \mathbf{y}}), & \mu_{\mathbf{x}}(0) &= \mathbf{1}_{\mathbf{x}=\mathbf{x}_0} \\ \gamma_{\mathbf{x}, \mathbf{y}} &= \mu_{\mathbf{x}} g_{\mathbf{x}, \mathbf{y}} \frac{\rho_{\mathbf{y}}}{\rho_{\mathbf{x}}}, & \rho_{\mathbf{x}} \neq 0, \mathbf{y} \neq \mathbf{x}. \end{aligned}$$

■

Note that this set of equations is identical to (3), but replaces the constant rate matrix  $\mathbb{Q}$  by a time varying matrix function  $\mathbb{G}(t)$ . Using this terminology, the first part of the fixed-point equations is

$$\eta^\alpha = \mathcal{R}(\mathbb{G}^\alpha, \mathbf{e}_0|_\alpha, \mathbf{e}_T|_\alpha), \quad (7)$$

where  $\mathbb{G}^\alpha(t)$  is the time-dependent matrix with entries

$$g_{\mathbf{x}_\alpha, \mathbf{y}_\alpha}^\alpha = \begin{cases} (q_{x_i y_i | \mathbf{u}_i}^{i | \mathbf{P} \mathbf{a}_i}) \mathbf{1}_{i \in A(\alpha)} \cdot n_{x_i, y_i}^{i \rightarrow \alpha} & \delta_{\mathbf{x}_\alpha, \mathbf{y}_\alpha} = \{i\} \\ \sum_{i \in N(\alpha)} \left( \mathbf{1}_{i \in A(\alpha)} q_{x_i x_i | \mathbf{u}_i}^{i | \mathbf{P} \mathbf{a}_i} + n_{x_i, x_i}^{i \rightarrow \alpha} \right) & \mathbf{x}_\alpha = \mathbf{y}_\alpha \\ 0 & \text{otherwise,} \end{cases} \quad (8)$$

and  $n^{i, \alpha}$  are time-dependent functions that originate from the Lagrange multipliers that enforce local consistency constraints,

$$\begin{aligned} \prod_{\alpha \in N(i)} n_{x_i, y_i}^{i \rightarrow \alpha} &= \left( \frac{\gamma_{x_i y_i}^i}{\mu_{x_i}^i} \right)^{c_i}, & x_i \neq y_i \\ \sum_{\alpha \in N(i)} n_{x_i, x_i}^{i \rightarrow \alpha} &= c_i \frac{\gamma_{x_i x_i}^i}{\mu_{x_i}^i}. \end{aligned} \quad (9)$$

These equations together with, (4) and (5) form the second set of equations that couple different clusters.

Equation (7) suggests that the matrix  $\mathbb{G}^\alpha$  plays the role of a rate matrix. Unlike  $\mathbb{Q}$ ,  $\mathbb{G}^\alpha$  is not guaranteed to be a rate matrix as its rows do not necessarily sum up to zero. Nonetheless, even though it is not a rate matrix, this system of equations has a unique solution that can be found using a backward-forward integration. Thus, since the matrix function  $\mathbb{G}^\alpha$  corresponds to a unique density set, we say that  $\mathbb{G}^\alpha$  is an *unnormalized parametrization* of the process  $P_{\eta^\alpha}$ .

At this point, it is tempting to proceed to construct a message passing algorithm based on this fixed point characterization. However, we are faced with a problem. Note that  $\lim_{t \rightarrow T} \frac{\gamma_{x_i e_i}}{\mu_{x_i}} = \infty$ . Therefore, according to Equation (9), when  $t$  approaches  $T$ , there exists some  $\alpha \in N(i)$  for which  $n_{x_i e_i}^{i, \alpha}(t)$  approaches  $\infty$  as  $t \rightarrow T$ . This implies that a simple-minded message passing procedure is susceptible to unbounded values and numerical difficulties.

### 4.3. Gauge Transformation

To overcome these numerical difficulties, we now derive an alternative characterization, which does not suffer from unbounded values. We start with a basic result.

**Proposition 4.4:** Let  $\mathbb{G}$  be a unnormalized rate matrix function, and let  $\omega_{\mathbf{x}}(t)$  be a smooth positive vector-valued function, where  $\omega_{\mathbf{x}}(t) > 0$  in  $[0, T)$ . Let  $\mathbb{G}^\omega$  to be the matrix function with

$$g_{\mathbf{x}\mathbf{y}}^\omega = \begin{cases} g_{\mathbf{x}\mathbf{y}} \cdot \frac{\omega_{\mathbf{x}}}{\omega_{\mathbf{y}}} & \mathbf{y} \neq \mathbf{x} \\ g_{\mathbf{x}\mathbf{x}} - \frac{d}{dt} \log \omega_{\mathbf{x}} & \mathbf{y} = \mathbf{x}. \end{cases} \quad (10)$$

Then,  $\mathcal{R}(\mathbb{G}, \mathbf{x}_0, \mathbf{x}_T) = \mathcal{R}(\mathbb{G}^\omega, \mathbf{x}_0, \mathbf{x}_T)$ .

**Proof sketch:** Let  $\rho, \eta$  satisfy the system of equations of Def. 4.3 with  $\mathbb{G}$ . Define  $\rho^\omega = \rho \cdot \omega$ , and show that  $\rho^\omega, \eta$  satisfy the same system of equations with  $\mathbb{G}^\omega$ . ■

This result characterizes transformations of (8–9) that do not change the fixed point solutions for cluster density sets. We seek transformations that reweigh the functions  $n^{i, \alpha}$  so that they remain bounded using the following result.

**Proposition 4.5:** Assume  $\mathbb{G}$  is a unnormalized rate matrix function such that  $g_{\mathbf{x}, \mathbf{y}}(t) \neq 0$  for all  $\mathbf{x}, \mathbf{y}$ ,  $g_{\mathbf{x}, \mathbf{y}}(t)$  is continuously differentiable in  $[0, T]$ , and  $\eta = \mathcal{R}(\mathbb{G}, \mathbf{x}_0, \mathbf{x}_T)$ . If  $\omega(t)$  is a family of smooth functions satisfying  $\omega_{\mathbf{x}}(T) = \mathbf{1}_{\mathbf{x}=\mathbf{x}_T}$  and  $\frac{d}{dt} \omega_{\mathbf{x}}(T) < 0$  for  $\mathbf{x} \neq \mathbf{x}_T$ , then

$$\lim_{t \rightarrow T} \frac{\gamma_{\mathbf{x}, \mathbf{y}}(t) \omega_{\mathbf{x}}(t)}{\mu_{\mathbf{x}}(t) \omega_{\mathbf{y}}(t)} < \infty, \quad \forall \mathbf{x} \neq \mathbf{y}$$

and

$$\lim_{t \rightarrow T} \left( \frac{\gamma_{\mathbf{x}, \mathbf{x}}(t)}{\mu_{\mathbf{x}}(t)} - \frac{d}{dt} \log \omega_{\mathbf{x}}(t) \right) < \infty, \quad \forall \mathbf{x}.$$

**Example 4.6:** One function that satisfies the conditions of Proposition 4.5 is  $\omega_{\mathbf{x}}(t) = 1 - t/T, \forall \mathbf{x} \neq \mathbf{e}_T$  and  $\omega_{\mathbf{e}_T}(t) = 1$ . ■

Using this result, we introduce weight functions  $\omega_{x_i}^i$  (as above) and define  $\omega_{\mathbf{x}_\alpha}^\alpha = \prod_{i \in N(\alpha)} (\omega_{x_i}^i)^{c_i / (c_i + 1)}$ . Using these weight functions, we define  $m_{x_i, y_i}^{i \rightarrow \alpha} = n_{x_i, y_i}^{i \rightarrow \alpha} \left( \frac{\omega_{x_i}^i}{\omega_{y_i}^i} \right)^{c_i / (c_i + 1)}$  and  $m_{x_i, x_i}^{i \rightarrow \alpha} = n_{x_i, x_i}^{i \rightarrow \alpha} - \frac{c_i}{c_i + 1} \frac{d}{dt} \log \omega_{x_i}^i$ .

**Algorithm 1** Continuous-Time Belief Propagation

---

Initialize messages: for all  $\alpha$  and all  $i \in N(\alpha)$   
 Choose  $\eta^\alpha \in \mathcal{M}_e^\alpha$   
 Compute  $\eta^{\alpha \rightarrow i}$  using (14)  
 Set  $m_{x_i, y_i}^{i \rightarrow \alpha} = 1 \forall x_i \neq y_i, m_{x_i, x_i}^{i \rightarrow \alpha} = 0$   
**repeat**  
 Choose a cluster  $\alpha$ :  
 1.  $\forall i \in N(\alpha)$ , set  $m_{x_i, y_i}^{i \rightarrow \alpha}$  using (15)  
 2. Update  $\tilde{\mathbb{G}}^\alpha$  using (11)  
 3. Compute  $\eta^\alpha$  from  $\tilde{\mathbb{G}}^\alpha$  using (12)  
 4.  $\forall i \in N(\alpha)$  compute  $\eta^{\alpha \rightarrow i}$  using (14)  
**until** convergence

---

Now if we define the time-dependent matrix  $\tilde{\mathbb{G}}^\alpha$  with entries

$$\tilde{g}_{\mathbf{x}_\alpha, \mathbf{y}_\alpha}^\alpha = \begin{cases} (q_{x_i y_i | \mathbf{u}_i}^i)^{\mathbf{I}_{i \in A(\alpha)}} \cdot m_{x_i, y_i}^{i \rightarrow \alpha} & \delta_{\mathbf{x}_\alpha, \mathbf{y}_\alpha} = \{i\} \\ \sum_{i \in N(\alpha)} \left( \mathbf{I}_{i \in A(\alpha)} q_{x_i x_i | \mathbf{u}_i}^i + m_{x_i, x_i}^{i \rightarrow \alpha} \right) & \mathbf{x}_\alpha = \mathbf{y}_\alpha \\ 0 & \text{otherwise,} \end{cases} \quad (11)$$

then  $\tilde{\mathbb{G}}^\alpha = (\mathbb{G}^\alpha)^{\omega^\alpha}$ . By Proposition 4.4,

$$\eta^\alpha = \mathcal{R}(\tilde{\mathbb{G}}^\alpha, \mathbf{e}_0 |_\alpha, \mathbf{e}_T |_\alpha). \quad (12)$$

Plugging the definition of  $m_{x_i, y_i}^{i \rightarrow \alpha}$  and  $m_{x_i, x_i}^{i \rightarrow \alpha}$  into (9) we get

$$\begin{aligned} \prod_{\alpha \in N(i)} m_{x_i, y_i}^{i \rightarrow \alpha} &= \left( \frac{\gamma_{x_i y_i}^i \omega_{x_i}^i}{\mu_{x_i}^i \omega_{y_i}^i} \right)^{c_i}, \quad x_i \neq y_i \\ \sum_{\alpha \in N(i)} m_{x_i, x_i}^{i \rightarrow \alpha} &= c_i \left( \frac{\gamma_{x_i x_i}^i}{\mu_{x_i}^i} - \frac{d}{dt} \log \omega_{x_i}^i \right). \end{aligned} \quad (13)$$

If the preconditions of Proposition 4.5 are satisfied, the terms in (13) are bounded. Together (11)–(13) provide an alternative characterization of the fixed point(s) of the optimization problem.

#### 4.4. Message Passing Scheme

We now use the above characterization as justification for a message passing scheme, that if converged, will satisfy the fixed point equations. While (11) and (12) are readily transformed into assignments, (13) poses a challenge.

We start by noting that (13) contains the terms  $\mu_{x_i}^i$  and  $\gamma_{x_i, y_i}^i$ . We can get these terms from  $\eta^\alpha$  for any  $\alpha \in N(i)$ . Thus, for  $\alpha \in N(i)$ , we define

$$\mu^{\alpha \rightarrow i} = \mu^\alpha |_i \quad \gamma^{\alpha \rightarrow i} = \gamma^\alpha |_i \quad (14)$$

We view these as the messages from cluster  $\alpha$  to the component  $i$ . At convergence,  $\mu^{\alpha \rightarrow i} = \mu^{\beta \rightarrow i}$  for  $\alpha, \beta \in N(i)$ , but this is not true before convergence.

Next, we rewrite (13) as an assignment

$$m_{x_i, y_i}^{i \rightarrow \alpha} = \begin{cases} \prod_{\substack{\beta \in N(i) \\ \beta \neq \alpha}} \frac{1}{m_{x_i, y_i}^{i \rightarrow \beta}} \frac{\gamma_{x_i y_i}^{\beta \rightarrow i} \omega_{x_i}^i}{\mu_{x_i}^i \omega_{y_i}^i} & x_i \neq y_i \\ \sum_{\substack{\beta \in N(i) \\ \beta \neq \alpha}} \left( \frac{\gamma_{x_i x_i}^{\beta \rightarrow i}}{\mu_{x_i}^i} - \frac{d}{dt} \log \omega_{x_i}^i - m_{x_i, x_i}^{i \rightarrow \beta} \right) & x_i = y_i \end{cases} \quad (15)$$

where we write  $\frac{\gamma_{x_i y_i}^i}{\mu_{x_i}^i} = \frac{\gamma_{x_i y_i}^{\beta \rightarrow i}}{\mu_{x_i}^i}$  once for each  $\beta$ .

The algorithm is summarized in Algorithm 1. The implementation of these steps involve a few details. We start with the initialization of messages. The only free parameter is the initial values of  $\eta^\alpha$ . To ensure that these initial choices are in  $\mathcal{M}_e^\alpha$ , we choose initial rates, and perform computations to get a valid posterior for the clusters. Another degree of freedom is the order of cluster updates. We use a randomized strategy, choosing a cluster at random, and if one of its neighbors was updated since it was last chosen, we update it.

The computation in Step 3, involves reverse integration followed by forward integration (as explained in Section 3). We gain efficiency by using adaptive numerical integration procedures. Specifically, we use the Runge-Kutta-Fehlberg (4,5) method (Press et al., 1992). This method chooses temporal evaluation points on the fly, and returns values at these points. The computations of Step 2 is done on demand only at the evaluation points. To allow efficient interpolation, we use a piecewise linear approximation of  $\eta$  whose boundary points are determined by the evaluation points that are chosen by the adaptive integrator. Finally, as might be expected, we do not have convergence guarantees. However, if the algorithm converges, the fixed point equations are satisfied, hence giving a stationary point (hopefully a local maximum) of problem (6).

## 5. Experiments

We tested our method on three representative network structures: a directed tree, a directed toroid, and a bidirectional ring (Fig. 2). The tree network does not have any cycles. The toroid network has cycles, but these are fairly long, whereas the bidirectional ring has multiple short cycles. All networks are parametrized as *dynamic Ising models*, in which neighboring components prefer to be in the same state. Specifically, we use the conditional rates

$$q_{x_i, y_i | \mathbf{u}_i}^i = \tau \left( 1 + \exp \left( -2y_i \beta \sum_{j \in \mathbf{P}\mathbf{a}_i} x_j \right) \right)^{-1}$$

where  $x_j \in \{-1, 1\}$ ,  $\beta$  is a *coupling parameter* and  $\tau$  is a *rate parameter*. Small values of  $\beta$  correspond to weak coupling, whereas  $\tau$  determines how fast components tend to switch states. For each experiment we set evidence at times 0 and 1 (Fig. 2, left panel).

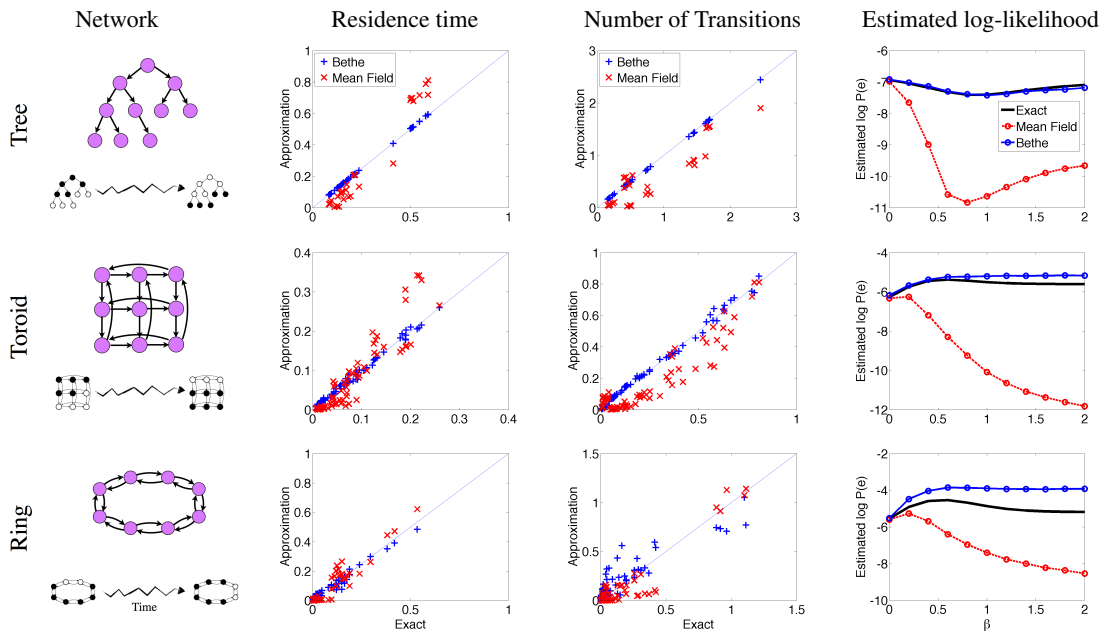


Figure 2. Simulation results for a tree network (top row), a toroid network (middle), and a bidirectional chain (bottom). **Left** network structure and the evidence at start and end points; black is +1 and white is -1. **Middle-left:** scatter plot of expected conditional residence times for networks with  $\beta = 1$ ,  $\tau = 8$ . Each point corresponds to a single statistic, the  $x$ -axis is the exact value and the  $y$ -axis is the approximate value. **Middle-right:** same for expected conditional transition times. **Right:** exact and approximations of log-likelihood as function of  $\beta$ , the strength of coupling between components ( $\tau = 8$ ).

We compare the Bethe approximation to exact inference and mean-field (Cohn et al., 2009). We start by comparing the value of sufficient statistics (residence time and number of transitions of each component for each state of its parents) computed by each method. For example, for a particular choice of  $\beta$  and  $\tau$ , (Fig. 2 middle columns) we can see that the Bethe approximation is virtually exact on the tree model and the toroid, but has some bias on the bidirectional ring model. These scatter plots also shed light on the nature of the difference between the two methods. Specifically, in the most likely scenario, two components switch from -1 to 1 near the beginning and the other two switch from 1 to -1 near the end, and so through most of the interval all the components are in the same state. The mean-field algorithm gives a uni-modal solution, focusing on the most likely scenario, resulting in zero residence time for the less likely states. These states are represented by the points close on the  $x$ -axis. The Bethe representation on the other hand can capture multiple scenarios.

Another aspect of the approximation is the estimation of the likelihood. In Fig. 2 (right column) we compare these estimations as function of  $\beta$ , the problem hardness. Again, we see that the Bethe approximation is essentially exact on the tree network, and provides close approximations in the two other networks. When we push  $\beta$  and  $\tau$  to extreme values we do see inaccuracies even in the tree network, showing

that the algorithm is an approximation.

While the ODE solvers used here allow adaptive integration error control, we do not have an a-priori control on the propagation of this error. To test this effect on overall accuracy, we repeated these experiments using standard grid refinement. Specifically, we computed integrals using uniformly spaced evaluation points and systematically halving integration intervals until no changes in the output were observed. Final results of these tests were practically the same as those obtained using adaptive integration.

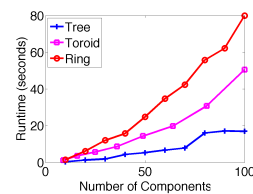


Figure 3. Run time vs. the number components in the three networks types ( $\beta = 1$ ,  $\tau = 8$ ).

Next, we examine how the algorithm scales with the number of components in the networks. In all three networks we see that the magnitude of relative error is essentially independent of the number of components (not shown). Fig. 3 shows that the run time scales linearly with the number of components. In harder networks the algorithm re-

quires more iterations leading to slower convergence.

## 6. Discussion

Here, we introduce a message passing scheme that provides a variational approximation for CTBNs. This scheme is derived from an approximate free energy functional, based on the principles of expectation-propagation, where we require the posteriors on clusters to be locally consistent in terms of the Markovian projections on individual components. We show that stationary points of this optimization problem are the fixed points of a message passing algorithm, whose structure closely mirrors Belief Propagation.

In contrast to Belief Propagation on discrete (atemporal) networks, our algorithm is *not* guaranteed to be exact on tree CTBNs. The source of inaccuracy is the projection of the marginal distributions over components onto Markov processes. While this projection loses information, our empirical results suggest that this approximation is relatively accurate. In contrast to methods in Dynamic Bayesian Networks (DBNs) (Koller & Friedman, 2009) that approximate the distribution over all components in each time slice, our approach approximates the temporal behavior of a single component over the whole time interval.

The works that are closest to ours are those of Nodelman et al. (2005) and Saria et al. (2007) which are also derived from an expectation-propagation energy functional. The main difference between the two methods is the structure of the approximation. Nodelman et al use a piecewise homogeneous representation, allowing them to represent the rate in each homogeneous segment by a constant (conditional) rate matrix. This, however, requires introducing machinery for deciding how to segment each component. As Saria et al show, this choice can have dramatic impact on the quality of the approximation and the running time. In contrast, our approach uses a (continuously) inhomogeneous representation, which is essentially the limit when segment sizes tend to zero. Surprisingly, rather than making the problem more complex, this choice simplifies the mathematics and also the implementation. In particular, our solution decouples the probabilistic issues (dependencies between components) and numerical issues (adaptive integration) and allows us to build on well-understood methods from numerical integration for efficient and adaptive selection of the number and placement of discretization points.

Our results show how a careful choice of representations and operations over them can narrow the gap between inference methods in discrete and continuous-time graphical models. Our constructions can be naturally generalized to capture more complex dependencies using methods based on Generalized Belief Propagation (Yedidia et al., 2005).

## Acknowledgments

We thank the anonymous reviewers for helpful remarks. This research was supported in part by a grant from the Israel Science Foundation. Tal El-Hay is supported by the Eshkol fellowship from the Israeli Ministry of Science.

## References

- Chung, K.L. *Markov chains with stationary transition probabilities*. 1960.
- Cohn, I., El-Hay, T., Friedman, N., and Kupferman, R. Mean field variational approximation for continuous-time Bayesian networks. UAI, 2009.
- El-Hay, T., Friedman, N., and Kupferman, R. Gibbs sampling in factorized continuous-time markov processes. UAI, 2008.
- Fan, Y. and Shelton, C. R. Learning continuous-time social network dynamics. UAI, 2009.
- Fan, Y. and Shelton, C.R. Sampling for approximate inference in continuous time Bayesian networks. In *AI & Math*, 2008.
- Felsenstein, J. *Inferring Phylogenies*. Sinauer, 2004.
- Gelfand, I. M. and Fomin, S. V. *Calculus of variations*. 1963.
- Heskes, T. and Zoeter, O. Expectation propagation for approximate inference in dynamic Bayesian networks. UAI, 2002.
- Jordan, M.I., Ghahramani, Z., Jaakkola, T.S., and Saul, L. K. An introduction to variational approximations methods for graphical models. In *Learning in Graphics Models*. 1998.
- Koller, D. and Friedman, N. *Probabilistic Graphical Models: Principles and Techniques*. 2009.
- Minka, T. P. Expectation propagation for approximate Bayesian inference. UAI, 2001.
- Nodelman, U., Shelton, C. R., and Koller, D. Continuous time Bayesian networks. UAI, 2002.
- Nodelman, U., Shelton, C.R., and Koller, D. Expectation propagation for continuous time Bayesian networks. UAI, 2005.
- Opper, M. and Sanguinetti, G. Variational inference for Markov jump processes. NIPS, 2007.
- Press, William H., Teukolsky, Saul A., Vetterling, William T., and Flannery, Brian P. *Numerical Recipes in C*. 1992.
- Saria, S., Nodelman, U., and Koller, D. Reasoning at the right time granularity. UAI, 2007.
- Simma, A., Goldszmidt, M., MacCormick, M., Barham, M., Black, M., Isaacs, M., and Mortier, R. CT-NOR: Representing and reasoning about events in continuous time. UAI, 2008.
- Wainwright, M. J. and Jordan, M. Graphical models, exponential families, and variational inference. *Found. Trends Mach. Learn.*, 1:1–305, 2008.
- Xu, J. and Shelton, C. R. Continuous time Bayesian networks for host level network intrusion detection. ECML/PKDD, 2008.
- Yedidia, J.S., Freeman, W.T., and Weiss, Y. Constructing free-energy approximations and generalized belief propagation algorithms. *IEEE Trans. Info. Theory*, 51:2282–2312, 2005.