

# A General Framework for Amortizing Variational Filtering

Joseph Marino, Milan Cvitkovic, Yisong Yue

California Institute of Technology (Caltech)

Caltech

## Problem

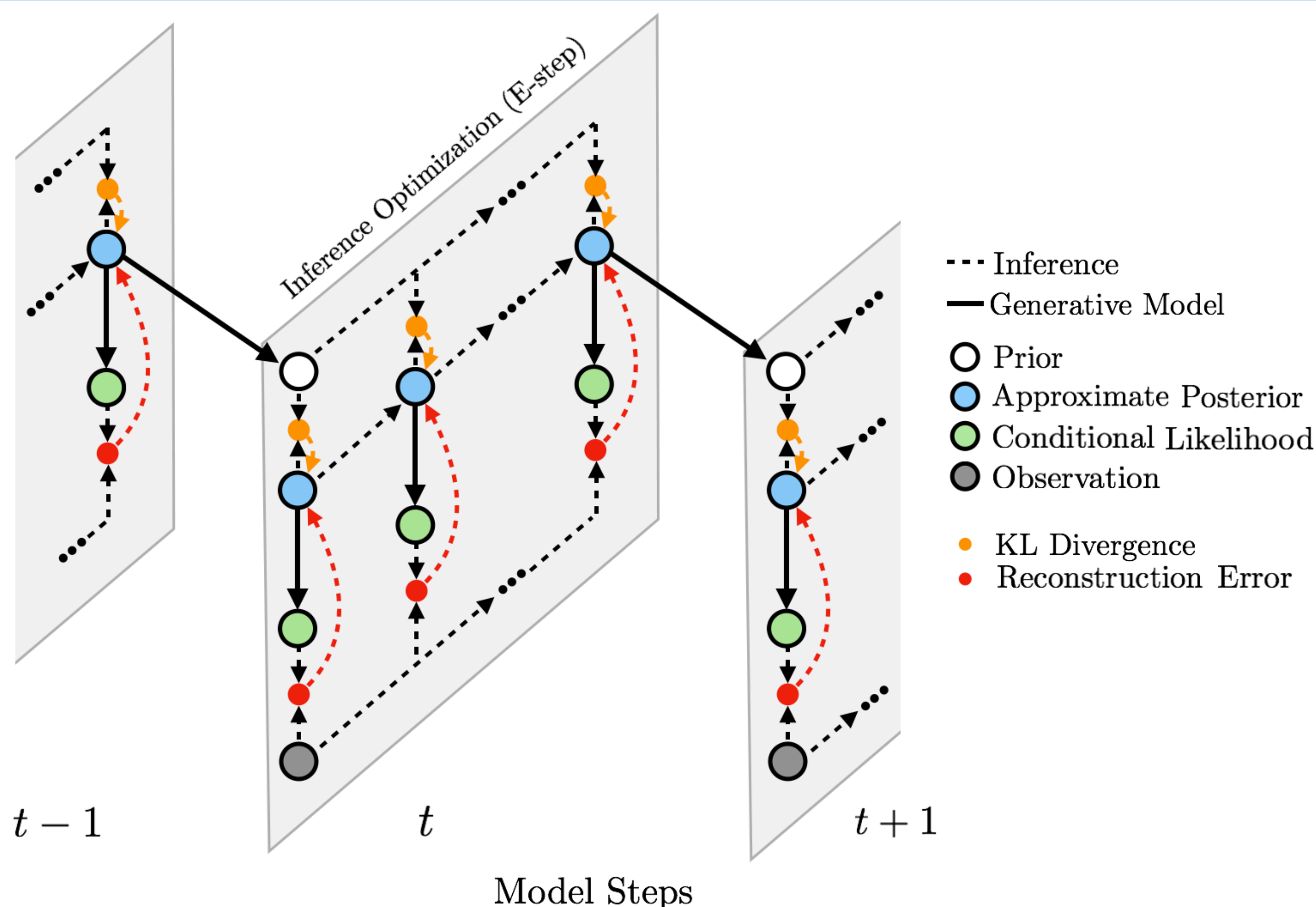
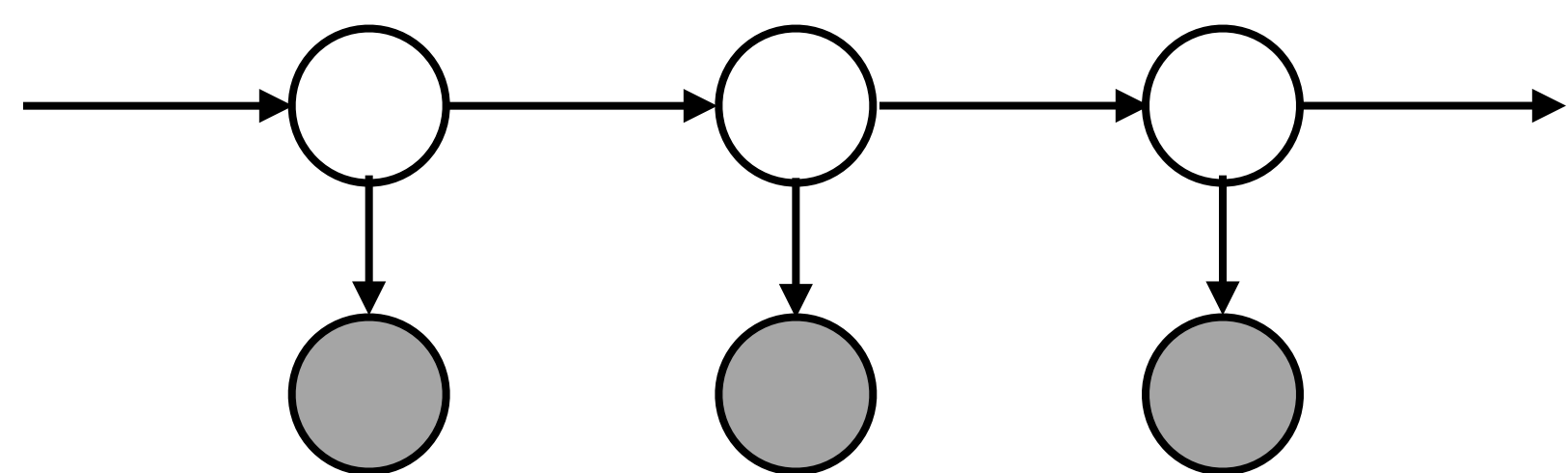
Deep latent variable models are often used for dynamical tasks, like reinforcement learning or time-series prediction. A central challenge is performing efficient online inference of the hidden states (filtering). In the static setting, amortized variational techniques are widely used for inference, but applying these techniques to dynamical problems has required hand-crafting an inference procedure for every new model. **We propose a general purpose method for efficiently performing accurate inference in any dynamical latent variable model.**

## Background

A dynamical latent variable model models a sequence of observations,  $\mathbf{x}_{\leq T}$ , using a sequence of latent variables,  $\mathbf{z}_{\leq T}$ , and parameters,  $\theta$ . These models are of the general form:

$$p_{\theta}(\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T}) = \prod_{t=1}^T p_{\theta}(\mathbf{x}_t | \mathbf{x}_{<t}, \mathbf{z}_{\leq t}) p_{\theta}(\mathbf{z}_t | \mathbf{x}_{<t}, \mathbf{z}_{<t}).$$

$p_{\theta}(\mathbf{x}_t | \mathbf{x}_{<t}, \mathbf{z}_{\leq t})$  is the *observation model*, and  $p_{\theta}(\mathbf{z}_t | \mathbf{x}_{<t}, \mathbf{z}_{<t})$  is the *dynamics model*. A simplified version of such models can be represented graphically as:



## Variational Filtering

Given a sequence of observations, we want to infer the posterior distribution over the sequence of latent variables,  $p_{\theta}(\mathbf{z}_{\leq T} | \mathbf{x}_{\leq T})$ . Unfortunately, this is often intractable. Instead, we use an approximate posterior,  $q(\mathbf{z}_{\leq T} | \mathbf{x}_{\leq T})$ , and minimize the following variational objective, called the **free energy**:

$$\mathcal{F} \equiv -\mathbb{E}_{q(\mathbf{z}_{\leq T} | \mathbf{x}_{\leq T})} \left[ \log \frac{p_{\theta}(\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T})}{q(\mathbf{z}_{\leq T} | \mathbf{x}_{\leq T})} \right].$$

We assume the **filtering** setting, where only past and present variables are used for inference, and assume the approximate posterior factorizes as

$$q(\mathbf{z}_{\leq T} | \mathbf{x}_{\leq T}) = \prod_{t=1}^T q(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{<t}).$$

With this filtering approximate posterior, the free energy becomes:

$$\mathcal{F} = \sum_{t=1}^T \mathbb{E}_{\prod_{\tau=1}^{t-1} q(\mathbf{z}_{\tau} | \mathbf{x}_{\leq \tau}, \mathbf{z}_{<\tau})} [\mathcal{F}_t]$$

$$\mathcal{F}_t \equiv -\mathbb{E}_{q(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{<t})} \left[ \log \frac{p_{\theta}(\mathbf{x}_t, \mathbf{z}_t | \mathbf{x}_{<t}, \mathbf{z}_{<t})}{q(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{<t})} \right]$$

### Algorithm 1 Variational Filtering Expectation Maximization

```

1: Input: observation sequence  $\mathbf{x}_{1:T}$ , model  $p_{\theta}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T})$ 
2:  $\nabla_{\theta} \mathcal{F} = 0$ 
3: for  $t = 1$  to  $T$  do
4:   initialize  $q(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{<t})$  ▷ from  $p_{\theta}(\mathbf{z}_t | \mathbf{x}_{<t}, \mathbf{z}_{<t})$ 
5:    $\tilde{\mathcal{F}}_t := \mathbb{E}_{q(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{<t})} [\mathcal{F}_t]$ 
6:    $q(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{<t}) = \arg \min_q \tilde{\mathcal{F}}_t$  ▷ inference (E-step)
7:    $\nabla_{\theta} \mathcal{F} = \nabla_{\theta} \mathcal{F} + \nabla_{\theta} \tilde{\mathcal{F}}_t$ 
8: end for
9:  $\theta = \theta - \alpha \nabla_{\theta} \mathcal{F}$  ▷ learning (M-step)

```

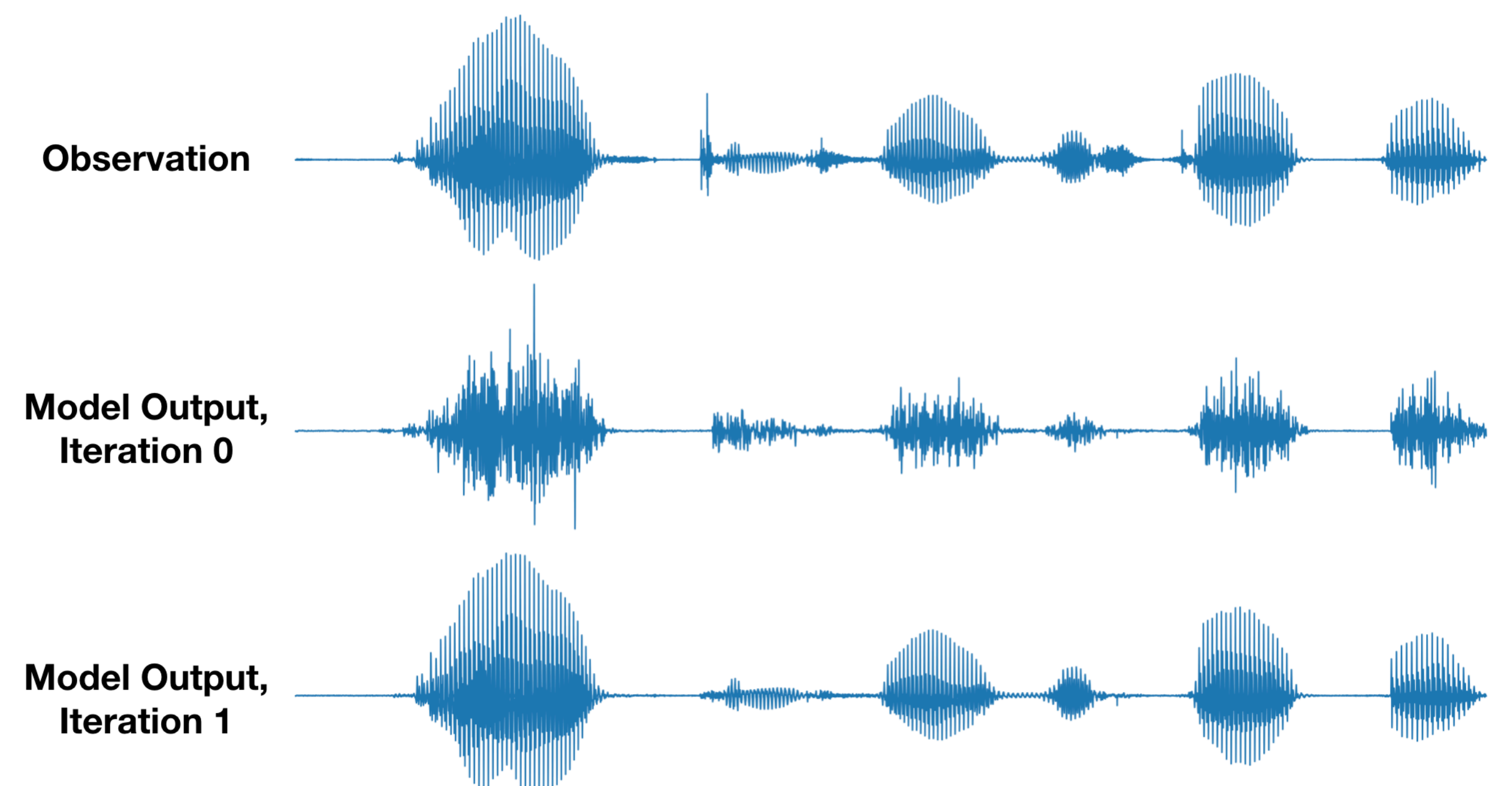
The **variational filtering EM** algorithm minimizes the filtering free energy by sequentially minimizing the free energy at each step. Initializing the approximate posterior at each step from the prior yields a Bayesian **prediction-update** loop.

We can amortize the inference optimization at each step by using an iterative inference model, which we refer to as **amortized variational filtering**. Denoting the approximate posterior parameters at step  $t$  as  $\lambda_t^q$ , the inference update is

$$\lambda_t^q \leftarrow f_{\phi}(\lambda_t^q, \nabla_{\lambda_t^q} \tilde{\mathcal{F}}_t).$$

## Results

We evaluate amortized variational filtering (**AVF**) on three dynamical latent variable models: **VRNN** (Chung *et al.*, 2015), **SRNN** (Fraccaro *et al.*, 2016), and **SVG** (Denton *et al.*, 2018). We train these models on speech, MIDI music, and video data sets.



We qualitatively demonstrate AVF on modeling the TIMIT speech data set. At the initial inference iteration of each step, the model outputs a *prediction* of the observation. At subsequent inference iterations, the output *reconstruction* is refined using approximate posterior gradients.

Music	Piano-midi.de	MuseData	JSB Chorales	Nottingham
SRNN				
baseline [Fraccaro et al., 2016]	8.20	6.28	4.74	2.94
baseline	8.19	6.27	6.92	3.19
AVF	<b>8.12</b>	<b>5.99</b>	<b>6.77</b>	<b>3.13</b>

Speech	TIMIT
VRNN	
baseline	1,082
AVF	<b>1,071</b>
SRNN	
baseline	1,026
AVF	<b>1,024</b>

Video	KTH Actions
SVG	
baseline	15,097
AVF	<b>11,714</b>

We compare AVF with baseline filtering methods for VRNN, SRNN, and SVG. We find that, in all cases, AVF results in improved model performance in terms of average free energy.

## Conclusion

Variational filtering EM and its amortized instantiation, amortized variational filtering, provide a **simple, theoretically motivated, and general-purpose method for performing filtering variational inference in dynamical latent variable models.**