

# Attend, Infer, Repeat: Notes

Adam Kosior

December 30, 2017

## 1 Introduction

Attend, Infer, Repeat (AIR; Eslami et al., 2016) is a VAE-type model capable of decomposing a static scene into its constituent parts. This is useful, since it explicitly represent distinct parts of the scene, which provides useful and actionable representations for down-stream processing. An example application are proposal location for object detection algorithms. While not guaranteed to represent objects, parts of the scene represented by separate representations are bound to be statistically significant. The purpose of this document is to explore the model in more detail.

## 2 Generative Model and Priors

Let  $\mathbf{x}$  an image,  $n$  a number of statistically interesting separate parts in the image,  $\mathbf{z} = \{\mathbf{z}^1, \dots, \mathbf{z}^n\}$  a group of variables, where each variable describes a part of the scene. Let  $p_\theta(\mathbf{z}, n) = p_\theta(\mathbf{z} | n)p_\theta(n)$  be a prior over latent variables describing our assumptions about possible layouts and complexity of the scene and let  $p_\theta(\mathbf{x} | \mathbf{z})$  be a *generating* model. The prior and the generating model together form a generative model of the scene, where the prior describes it in terms of latent variables and the generating model uses the latent variables to paint the scene. We can write down the marginal distribution of images as

$$p_\theta(\mathbf{x}) = \sum_{n=1}^N p_{\theta, N}(n) \int p_\theta(\mathbf{z} | n) p_\theta(\mathbf{x} | \mathbf{z}) d\mathbf{z}. \quad (1)$$

We take the prior to have the following form:

$$p_\theta(n) = \text{Geom}(n | \theta), \quad (2)$$

$$p_\theta(\mathbf{z} | n) = \prod_{i=1}^n p_\theta(\mathbf{z}^i) = \prod_{i=1}^n \mathcal{N}(\mathbf{o}, \mathbf{I}), \quad (3)$$

which assumes that  $\mathbf{z}^i$  are independent under the prior. The success probability for the geometric distribution is chosen to encourage sparse (as few steps as possible) solutions. In practice, it is annealed over 100k training steps to a very low value on the order of  $p = 10^{-e}$ ,  $e \in [5, 10]$ .

### 3 Inference through an Approximate Posterior

Since the latent variables  $\mathbf{z}^i$  are independent under the prior, they are exchangeable. It introduces symmetries, which increases the volume of the search space in a combinatorial way. To address this issue, the original work formulates inference as a sequential problem, where latent variables describing a part of a scene depend on previously inferred latent variables. Specifically, it parametrises  $\mathbf{n}$  as a vector of  $n$  ones followed by a zero  $\mathbf{z}^p$ . With this representation and a sequential implementation of inference, it is enough to output a single *presence* indicator variable at each processing step. The approximate posterior is given by

$$q_\phi(\mathbf{z}, \mathbf{z}^p | \mathbf{x}) = q_\phi(z^{p,n+1} = 0 | \mathbf{z}^{1:n}, \mathbf{x}) \prod_{i=1}^n q_\phi(\mathbf{z}^i, z^{p,i} = 1 | \mathbf{z}^{1:i-1}, \mathbf{x}), \quad (4)$$

with  $q_\phi$  implemented as a neural network. To avoid explaining the same object twice, it is vital to capture the dependency  $\mathbf{z}^i, z^{p,i} | \mathbf{z}^{1:i-1}, \mathbf{x}$ , which is implemented as a recurrent neural network  $R_\phi$  with hidden state  $\mathbf{h}^i$ , where the initial hidden state  $\mathbf{h}^0$  is randomly initialised and learnable, and specifically as

$$\omega^i, \mathbf{h}^i = R_\phi(\mathbf{x}, \mathbf{z}^{i-1}, \mathbf{h}^{i-1}). \quad (5)$$

The variable  $\omega^i$  specifies parameters of the probability distribution over  $\mathbf{z}^i$  and  $z^{p,i}$  and introduces a conditional independence property, namely  $\mathbf{z}^i \perp\!\!\!\perp z^{p,i} | \omega^i$ . We can use this property to factorise the approximate posterior distribution as

$$q_\phi(\mathbf{z}, \mathbf{z}^p | \mathbf{x}) = q_\phi(z^{p,n+1} = 0 | \omega^{n+1}) \prod_{i=1}^n q_\phi(\mathbf{z}^i | \omega^i) q_\phi(z^{p,i} | \omega^i). \quad (6)$$

Since the hidden state  $\mathbf{h}^i$  is deterministic, it follows a Dirac distribution centred on  $\mathbf{h}^i$  and there is no need to integrate over it: the integral collapses to a single value. Latent distribution statistics  $\omega$  are computed iteratively and depend on the previous latent variables. This precludes analytical computation of the KL-divergence, which has to be approximated by sampling.

## 4 Learning by Maximising the ELBO

AIR is trained by maximising the evidence lower bound  $\mathcal{L}_{\text{AIR}}$  (ELBO) given by

$$\mathcal{L}_{\text{AIR}}(\phi, \theta) = \mathbb{E}_{q_\phi} \left[ \log \frac{p_\theta(\mathbf{z}, \mathbf{n}, \mathbf{x})}{q_\phi(\mathbf{z}, \mathbf{n} | \mathbf{x})} \right] \quad (7)$$

$$= \mathbb{E}_{q_\phi} [\log p_\theta(\mathbf{x} | \mathbf{z})] - \text{KL}(q_\phi(\mathbf{z}, \mathbf{n} | \mathbf{x}) \| p_\theta(\mathbf{z}, \mathbf{n})). \quad (8)$$

The first term is a probabilistic analog of the reconstruction error, while the second term is a complexity penalty. In this case it encourages minimum-length encoding of every part of the scene as well as decomposing the scene into a minimal number of parts. The KL term can be rewritten according to the product rule as

$$\text{KL}(q_\phi(\mathbf{z}, \mathbf{n} | \mathbf{x}) \| p_\theta(\mathbf{z}, \mathbf{n})) = \quad (9)$$

$$= \text{KL}(q_\phi(\mathbf{n} | \mathbf{x}) \| p_\theta(\mathbf{n})) + \mathbb{E}_{q_\phi(\mathbf{n} | \mathbf{x})} [\text{KL}(q_\phi(\mathbf{z} | \mathbf{n}, \mathbf{x}) \| p_\theta(\mathbf{z} | \mathbf{n}))]. \quad (10)$$

## 5 Policy Gradient Techniques

Policy gradients are used here to backprop through samples from a discrete probability distribution.

### 5.1 Score Function

Williams, 1992

### 5.2 NVIL

Mnih and Gregor, 2014 implemented but doesn't work that well, presumably because of different effective weightings of the continuous and discrete components of the loss. I still have to check with the lower learning rate.

### 5.3 Discrete IWAE

To be done... Originally introduced by Burda, Grosse, and Salakhutdinov, 2015, but extended to a discrete case by Tucker et al., 2017.

### 5.4 Convex Relaxation

To be done... concrete distribution/gumbel Jang, Gu, and Poole, 2016; Madison, Mnih, and Teh, 2016

## 5.5 REBAR

This seems terribly complicated Tucker et al., 2017

## 6 FIVO for Better Bounds

Introduced by Maddison, Lawson, et al., 2017; ywt says its much better for sequential estimation such as AIR than the usual ELBO.

## 7 Complex Backgrounds: Behaviour and Extensions

AIR depends on statistical differences between pixels belonging to an object and the ones belonging to the background. Moreover, the way it has been defined, it can reconstruct only objects while neglecting the background completely. It is natural to expect that it simply wouldn't work in the case of non-empty backgrounds. To mitigate this issue, it is necessary to augment the model with a background-explaining component  $g$ . In a simple case, it could produce a background  $\mathbf{b}$  based on the last hidden state of the controller RNN  $R_\phi$ , specifically:

$$\mathbf{b} = g_\phi(\mathbf{h}^N) \quad (11)$$

$$\mathbf{y} = h^{\text{dec}}(\mathbf{z}, \mathbf{b}) \quad (12)$$

For the case of MNIST,  $g$  can be implemented as a simple MLP. For more complicated model, it would have to be a more complicated model based on up-convolutions, but this is also true about the glimpse decoder.

Design an experiment to test this behaviour.

## 8 An Extension to Timeseries

AIR displays an inconsistent behaviour depending on initialisation. It is randomly initialised and it tends to converge to different solutions when training is restarted. Only a small subset of the solutions is close to what the user would want, namely explaining one object at a time. This is understandable, since a dataset of single images cannot encode an *consistency* prior, whereby objects would have to be preserved when they move. Extending AIR to sequential data has a potential of solving this issue, since it should be much easier to explain an object by changing its position parameter than to inferring its appearance description from scratch at every time-step. Moreover, it prevents explaining two objects moving separately with a single glimpse, since it would be impossible to preserve them consistently while reconstructing.

An extension to timeseries is not trivial, however. It requires some form of dependency between consecutive timestep. It is important to distinguish temporal and within-timestep dependencies. We dub the latter *sequential* dependencies. In the original work, the only sequential dependencies exist between the hidden states  $\mathbf{h}^i \mid \mathbf{h}^{i-1}$  and the presence variables  $z^{p,i} \mid z^{p,1:i-1}$ , while there is no sequential dependence between appearance and position  $\mathbf{z}^i$ . A temporal version of the algorithm requires (or might benefit from) the following:

- between-timestep hidden state dependence e.g.,  $\mathbf{h}^t \mid \mathbf{h}^{t-1}$ , where this timestep might or might not be the same as the sequential hidden state
- temporal dependence on object position
- temporal dependence on object appearance
- temporal dependence on object presence

### 8.1 Indian Buffer Process for Object Counts

Read and write (Gael, Teh, and Ghahramani, 2009).

## 9 Sequential AIR

Let us start by describing the model from the point of view of its generative story. Images are generated by first assuming that, at every time-step, objects are propagated from previous time-step and some new objects can be introduced. Let  $k \in \{0, 1, \dots, K\}$ , the number of objects propagated from the previous time-step and let  $n \in \{0, 1, \dots, N - k\}$  the number of objects discovered at the current time-step. The maximum of  $K \in \mathcal{N}_+$  objects can be propagated and the model can handle up to  $N \in \mathcal{N}_+$  total objects, therefore  $K \leq N$ . Moreover, at every time-step, the model will generate a set of latent variables, one for each object. If we use superscripts D and P to denote latent variables generated by discovery and propagation models, respectively, then the aggregated latent variable  $\mathbf{z}_t = \mathbf{z}_t^P \cup \mathbf{z}_t^D$ , where  $\mathbf{z}_t^P = \{\mathbf{z}_t^{P,1}, \dots, \mathbf{z}_t^{P,k_t}\}$  and  $\mathbf{z}_t^D = \{\mathbf{z}_t^{D,1}, \dots, \mathbf{z}_t^{D,k_t}\}$ . Consequently, if we use  $m_t = k_t + n_t$  to denote the total number of objects present at time  $t$ , then  $\mathbf{z}_t = \{\mathbf{z}_t^1, \dots, \mathbf{z}_t^{m_t}\}$ . We note that the structure of the latent variables from the discovery and propagation models is the same, and in its simplest form it consists of  $\mathbf{z}_t^i = \{\mathbf{z}_t^{\text{what},i}, \mathbf{z}_t^{\text{where},i}, b_t^j\}$ . The components represent appearance, location and presence of an object at the given time-step.

At the first time-step  $t = 1$  there are no objects to propagate, so we sample up to  $N$  objects from a discovery prior  $p(n_1, \mathbf{z}_1^D \mid N)$ . Starting from the second time-step  $t = 2$ , the model first propagates objects by sampling from

a propagation prior  $p(k_2, \mathbf{z}_2^P | n_1 + k_1, \mathbf{z}_1)$ , where  $\mathbf{z}_1 = \mathbf{z}_1^D$ . The model also samples  $n_2$  new objects from the prior  $p(n_2, \mathbf{z}_2^D | N - k_2)$ . From now on, that is for  $t \geq 2$ , we set the aggregated latent variable  $\mathbf{z}_t = \{\mathbf{z}_t^P, \mathbf{z}_t^D\}$ . This process continues up to the final time-step  $T$ . The images are generated by passing the latent variables  $\mathbf{z}_{1:T}$  to the generating model  $p_\theta(\mathbf{x}_t | \mathbf{z}_t)$  one at a time. Note that  $\mathbf{z}_t$  is a set of up to  $N$  latent variables, where each latent variable represents a separate object. The generating model acts separately on every latent variables in the set, and the output random variable of the generating model consists of the sum of outputs of the generative model for each latent variable in the set. Figure 1 shows the graphical model of the generative story. Please note that the prior for the number of new objects stays the same for every time-step and it is analogous to the prior used by AIR, with the exception that the number of maximum objects can differ.

To make it more formal, we will now give all equations governing the generative process. The prior for the discovery model is defined as follows.

$$p(n_t, \mathbf{z}_t^D | N) = p(n_t | N) p(\mathbf{z}_t | n_t) = \text{TruncatedGeom}(n_t | N) \prod_{i=1}^{n_t} \mathcal{N}(\mathbf{o}, \mathbf{I}), \quad (13)$$

where  $\text{TruncatedGeom}(n_t | N)$  is a geometric distribution, whose support is truncated to  $\{0, 1, \dots, N\}$ . The propagation model is defined by the following.

$$p(k_t, \mathbf{z}_t^P | m_{t-1}, \mathbf{z}_{1:t-1}) = \prod_{j=1}^m p(b_t^{P,j} | \mathbf{z}_{1:t-1}^j) p(\mathbf{z}_t^j | \mathbf{z}_{1:t-1}^j) = \prod_{j=1}^{m_{t-1}} \text{Bernoulli}(b_t^{P,j} | \psi(\mathbf{z}_{1:t-1}^j)) \mathcal{N}(\mathbf{z}_t^j | \boldsymbol{\mu}(\mathbf{z}_{1:t-1}^j), \boldsymbol{\sigma}^2(\mathbf{z}_{1:t-1}^j)), \quad (14)$$

with  $m_t = n_t + k_t$  the total number of objects present at time  $t$ ,  $\mathbf{z}_{1:t-1}^j$  the history of latent variables belonging to object  $j$  and  $\mathbf{z}_{1,0}^j = \emptyset$ . In the above we reparameterise  $k_t$  as a sum of Bernoulli random variables  $b_t^{P,j}$ , that is  $k_t = \sum_{i=1}^m b_t^{P,i}$ . Since the probability of propagation  $\psi^j$  is different for every object,  $k_t$  follows a Poisson-Binomial distribution with  $m$  trials. Direct conditioning on the whole history of latent variables is impractical, and therefore we use a learned deterministic RNN  $R_\theta^{\text{prior}}$  to estimate  $\psi$ ,  $\mu$  and  $\sigma^2$ . the RNN shares parameters but maintains a separate hidden state for every object  $j$ . We can rewrite Equation (14) in terms of a simple recurrence as

$$p(k_t, \mathbf{z}_t^P | m_{t-1}, \mathbf{z}_{1:t-1}) = p(k_t, \mathbf{z}_t^P | m_{t-1}, \mathbf{h}_t^{\text{prior}, j}) = \prod_{j=1}^{m_{t-1}} \text{Bernoulli}(b_t^{P,j} | \psi(\mathbf{h}_t^{\text{prior}, j})) \mathcal{N}(\mathbf{z}_t^j | \boldsymbol{\mu}(\mathbf{h}_t^{\text{prior}, j}), \boldsymbol{\sigma}^2(\mathbf{h}_t^{\text{prior}, j})), \quad (15)$$

$$\mathbf{h}_t^{\text{prior}, j} = \mathbb{R}_\theta^{\text{prior}} \left( \mathbf{z}_{t-1}^j, \mathbf{h}_{t-1}^{\text{prior}, j} \right). \quad (16)$$

Any parameters in the generating model, e.g., the parameters of the RNN, can be learnt jointly with the inference model (to be described shortly). Alternatively, one could use a non-parametric prior for the discrete latent variable  $k_{1:T}$ , e.g., the Indian Buffet Process.

Images are created by decoding the latent variables into small ‘glimpses’ that are then transformed to match the size of the original image and summed up. Each latent variable  $\mathbf{z}_t^j$  can be decomposed as  $\mathbf{z}_t^j = \left\{ \mathbf{z}_t^{j, \text{where}}, \mathbf{z}_t^{j, \text{what}} \right\}$ . More formally,

$$\hat{\mathbf{g}}_t^j = f_\theta^{\text{dec}} \left( \mathbf{z}_t^{\text{what}, j} \right), \quad (17)$$

$$\mathbf{y}_t^j = f_\theta^{\text{STN}} \left( \hat{\mathbf{g}}_t^j, \mathbf{z}_t^{\text{where}, j} \right), \quad (18)$$

$$\boldsymbol{\mu}_x = \sum_{j=1}^{m_t} \mathbf{y}_t^j, \quad (19)$$

$$\hat{\mathbf{x}}_t \sim \mathcal{N} \left( \mathbf{x}_t \mid \boldsymbol{\mu}_x, \boldsymbol{\sigma}_x^2 \right), \quad (20)$$

with the decoder  $f_\theta^{\text{dec}}$ , spatial transformer  $f_\theta^{\text{STN}}$  and a fixed per-pixel standard deviation  $\boldsymbol{\sigma}_x^2$ .

All components of the generative model taken together specify a joint probability  $p_\theta(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}, \mathbf{n}_{1:T}, k_{1:T})$  over the sequence of images and the corresponding latent variables. We can expand it as

$$\begin{aligned} p_\theta(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}, \mathbf{n}_{1:T}, k_{1:T}) &= \prod_{t=1}^T p_\theta(\mathbf{x}_t, \mathbf{z}_t, \mathbf{n}_t, k_t \mid \mathbf{z}_{1:t-1}, \mathbf{n}_{1:t-1}, k_{1:t-1}) \\ &= \prod_{t=1}^T p_\theta(\mathbf{x}_t \mid \mathbf{z}_t) p_\theta(\mathbf{z}_t, \mathbf{n}_t, k_t \mid \mathbf{z}_{1:t-1}, \mathbf{n}_{1:t-1}, k_{1:t-1}) \\ &= \prod_{t=1}^T \underbrace{p_\theta(\mathbf{x}_t \mid \mathbf{z}_t)}_{\text{generation}} \underbrace{p_\theta(\mathbf{z}_t^D, \mathbf{n}_t \mid \mathbf{z}_t^P, k_t)}_{\text{discovery}} \underbrace{p_\theta(\mathbf{z}_t^P, k_t \mid \mathbf{z}_{1:t-1}, \mathbf{n}_{1:t-1}, k_{1:t-1})}_{\text{propagation}}. \end{aligned} \quad (21)$$

In the second line we used the fact that the observation at time  $t$  is conditionally independent of latent variables at previous time-steps given the current ones. In the third line we decomposed  $\mathbf{z}_t$  as the union of latent variables corresponding to the discovered objects  $\mathbf{z}_t^D$  and propagated objects  $\mathbf{z}_t^P$ . Only the latter depend on the history of latent variables; the former depend only on what is captured in propagation to the current step. The last term can be

decomposed as

$$\begin{aligned}
& p_\theta(\mathbf{z}_t^P, k_t \mid \mathbf{z}_{1:t-1}, n_{1:t-1}, k_{1:t-1}) \\
&= \int p_\theta(\mathbf{z}_t^P, k_t \mid \mathbf{z}_{t-1}, n_{t-1} + k_{t-1}, \mathbf{h}_t^P) p_\theta(\mathbf{h}_t^P \mid \mathbf{h}_{t-1}^P, \mathbf{z}_{t-1}) d\mathbf{h}_t^{\text{prior}} \\
&= p_\theta(\mathbf{z}_t^P, k_t \mid \mathbf{z}_{t-1}, m_{t-1}, \mathbf{h}_t^{\text{prior}}) \big|_{\mathbf{h}_t^{\text{prior}} = \mathbf{R}_\theta^{\text{prior}}(\mathbf{h}_{t-1}^{\text{prior}}, \mathbf{z}_{t-1})},
\end{aligned} \tag{22}$$

where we introduce the hidden state  $\mathbf{h}_t^P$  of the RNN used for propagating objects through time. Computation of the state is deterministic and therefore  $p_\theta(\mathbf{h}_t^{\text{prior}} \mid \mathbf{h}_{t-1}^{\text{prior}}, \mathbf{z}_{t-1})$  is a Dirac distribution centred on  $\mathbf{h}_t^P$ , which further simplifies computation. In the second line we also used the fact that object propagation depends only on the number of objects present at the previous time-step and it does not matter whether the objects were propagated from previous time-step or newly discovered.

## 9.1 Inference

To infer the latent variables, we have to invert the generative model. To do so, the inference model has to take into account the observations  $\mathbf{x}_{1:T}$  and produce sequences of sets of latent variables  $\mathbf{z}_{1:T}$ . Similar to the generative process, inference uses a two-step approach, whereby it propagates objects from previous time-steps and then tries to discover new objects. Let us start by defining the approximate posterior distribution over all latent variables.

$$\begin{aligned}
& q_\phi(\mathbf{z}_{1:T}, n_{1:T}, k_{1:T} \mid \mathbf{x}_{1:T}) \\
&= \prod_{t=1}^T q_\phi(\mathbf{z}_t, n_t, k_t \mid \mathbf{x}_{1:t}, \mathbf{z}_{1:t-1}, n_{1:t-1}, k_{1:t-1}) \\
&= \prod_{t=1}^T q_\phi(\mathbf{z}_t, n_t, k_t \mid \mathbf{x}_t, \mathbf{z}_{1:t-1}, n_{1:t-1}, k_{1:t-1}) \\
&= \prod_{t=1}^T q_\phi(\mathbf{z}_t, n_t, k_t \mid \mathbf{x}_t, \mathbf{z}_{t-1}, n_{t-1}, k_{t-1}, \mathbf{h}_t^q) \big|_{\mathbf{h}_t^q = \mathbf{f}_q^{\text{RNN}}(\mathbf{z}_{t-1}, \mathbf{h}_{t-1}^q)} \\
&= \prod_{t=1}^T q_\phi(\mathbf{z}_t^D, n_t \mid \mathbf{x}_t, \mathbf{z}_t^P, k_t) q_\phi(\mathbf{z}_t^P, k_t \mid \mathbf{x}_t, \mathbf{z}_{t-1}, \mathbf{k}_{t-1} + \mathbf{n}_{t-1}, \mathbf{h}_t^q) \\
&= \prod_{t=1}^T \underbrace{q_\phi^D(\mathbf{z}_t^D, n_t \mid \mathbf{x}_t, \mathbf{z}_t^P, k_t)}_{\text{discovery}} \underbrace{q_\phi^P(\mathbf{z}_t^P, k_t \mid \mathbf{x}_t, \mathbf{z}_{t-1}, \mathbf{m}_{t-1}, \mathbf{h}_t^q)}_{\text{propagation}}
\end{aligned} \tag{23}$$

In the second line we use the product rule to rewrite the joint distribution as a temporal chain. In the third line we introduce the assumption that latent



variables are conditionally independent of the past observations given the corresponding latent variables. In the fourth line we substituted the dependency on the history of latent variables by introducing a Markovian hidden state  $\mathbf{h}_t^q$ . We assume that the hidden state follows a Dirac distribution and therefore the associated integral collapses to a single value. In line five we decompose the latent variable  $\mathbf{z}_t$  into the ones introduced by the discovery and propagation steps, respectively. We also omit the value of the hidden state to reduce clutter. In line six we explicitly condition the propagation step on the number of latent variables in the previous time-step.

We start by describing the propagation model  $q_\phi$ . We would like to condition propagation of objects on all other objects, and therefore we decompose inference into a sequence of steps.

$$\begin{aligned}
q_\phi^p(\mathbf{z}_t^p, k_t \mid \mathbf{x}_t, \mathbf{z}_{t-1}, m_{t-1}, \mathbf{h}_t^q) \\
&= \prod_{i=1}^{m_{t-1}} q_\phi^p(\mathbf{z}_t^{p,i}, b_t^{p,i} \mid \mathbf{x}_t, \mathbf{z}_{t-1}^i, \mathbf{h}_t^q, \mathbf{z}_t^{p,1:i-1}) \\
&= \prod_{i=1}^{m_{t-1}} q_\phi^p(\mathbf{z}_t^{p,i}, b_t^{p,i} \mid \mathbf{x}_t, \mathbf{z}_{t-1}^i, \mathbf{h}_t^q, \mathbf{h}_t^{p,i}) \Big|_{\mathbf{h}_t^{p,i} = \mathbb{R}_\phi^p(\mathbf{z}_t^{p,i-1}, b_t^{p,i-1} \mathbf{h}_t^{p,i-1})}
\end{aligned} \tag{24}$$

$\mathbf{z}_t^{p,0}$  and  $b_t^{p,0}$  start as zeros and  $\mathbf{h}_t^{p,0}$  is randomly initialised and learnable. It is worth noting that  $m_t = \sum_{i=1}^{m_{t-1}} b_t^{p,i} \leq m_{t-1}$ , therefore propagation can forget objects, but it cannot introduce any new ones. At the end of the propagation stage, the hidden state  $\mathbf{h}_t^{S,k_t}$  should contain information about all the propagated objects and can be used to condition the discovery stage.

To handle object discovery, we follow the implementation of AIR, with the difference that we condition on the hidden state of the propagation phase  $\mathbf{h}_t^{p,k_t}$  and that there is a variable maximum number of steps. It defines a sequential inference procedure, where the model attends to and explains different regions of the image, one at a time. The approximate posterior distribution is defined by Equations (4) to (6). We redefine it here to account for temporal dependencies and conditioning. The approximate posterior distribution for

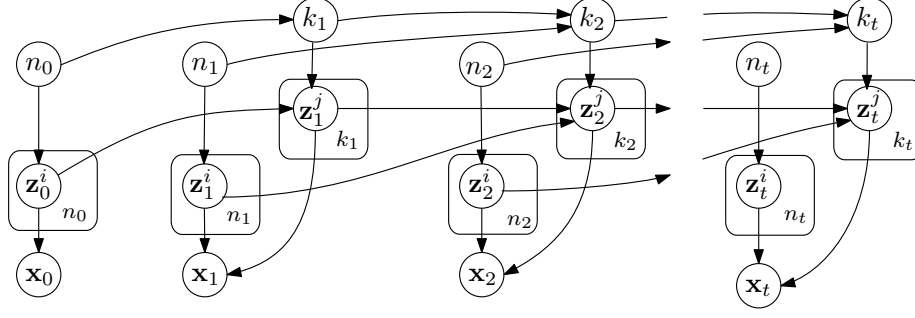


Figure 1: The generative story of sequential AIR.

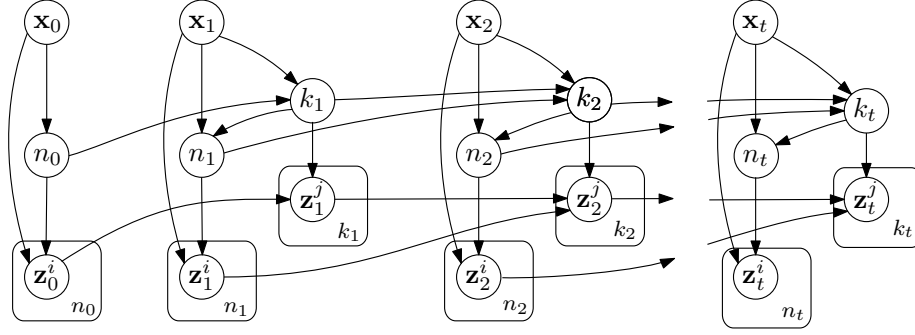


Figure 2: Graphical model for inference in sequential AIR.

discovered objects is given by

$$\begin{aligned}
& q_\phi^D(z_t^D, n_t \mid x_t, \mathbf{z}_t^P, k_t) \\
&= q_\phi^D(z_t^D, n_t \mid x_t, \mathbf{h}_t^{P, k_t}, k_t) \\
&= q_\phi^D(b_t^{D, n_t+1} = 0 \mid x_t, \mathbf{z}_t^{D, 1:n_t}, \mathbf{h}_t^{P, k_t}) \prod_{i=1}^{n_t} q_\phi^D(b_t^{D, i}, z_t^{D, i} \mid x_t, \mathbf{z}_t^{D, 1:i-1}, \mathbf{h}_t^{P, k_t}) \\
&= q_\phi^D(b_t^{D, n_t+1} = 0 \mid x_t, \mathbf{h}_t^{D, n_t}, \mathbf{h}_t^{P, k_t}) \prod_{i=1}^{n_t} q_\phi^D(b_t^{P, i}, z_t^{D, i} \mid x_t, \mathbf{h}_t^{D, i-1}, \mathbf{h}_t^{P, k_t}),
\end{aligned}$$

where  $\mathbf{h}_t^{D, i} = \mathbf{R}_\phi^D(z_t^{D, i}, \mathbf{h}_t^{D, i-1})$ .

(25)

## References

- Burda, Yuri, Roger Grosse, and Ruslan Salakhutdinov (2015). “Importance Weighted Autoencoders”. In: ISSN: 1312.6114v10. arXiv: 1509.00519. URL: <http://arxiv.org/abs/1509.00519>.
- Eslami, S. M. Ali et al. (2016). “Attend, Infer, Repeat: Fast Scene Understanding with Generative Models”. In: *NIPS*. arXiv: 1603.08575. URL: <http://arxiv.org/abs/1603.08575>.
- Gael, Jurgen Van, Yee Whye Teh, and Zoubin Ghahramani (2009). “The Infinite Factorial Hidden Markov Model”. In: *NIPS*, pp. 1697–1704. ISBN: 9781605609492. URL: <https://papers.nips.cc/paper/3518-the-infinite-factorial-hidden-markov-model>.
- Jang, Eric, Shixiang Gu, and Ben Poole (2016). “Categorical Reparameterization with Gumbel-Softmax”. In: ISSN: 1611.01144. arXiv: 1611.01144. URL: <http://arxiv.org/abs/1611.01144>.
- Maddison, Chris J., Dieterich Lawson, et al. (2017). “Filtering Variational Objectives”. In: *NIPS*. arXiv: 1705.09279. URL: <http://arxiv.org/abs/1705.09279>.
- Maddison, Chris J., Andriy Mnih, and Yee Whye Teh (2016). “The Concrete Distribution: A Continuous Relaxation of Discrete Random Variables”. In: arXiv: 1611.00712. URL: <http://arxiv.org/abs/1611.00712>.
- Mnih, Andriy and Karol Gregor (2014). “Neural Variational Inference and Learning in Belief Networks”. In: *ICML*. ISBN: 9781634393973. arXiv: arXiv: 1402.0030v2. URL: <http://arxiv.org/abs/1402.0030>.
- Tucker, George et al. (2017). “REBAR: Low-variance, unbiased gradient estimates for discrete latent variable models”. In: *NIPS*. arXiv: 1703.07370. URL: <http://arxiv.org/abs/1703.07370>.
- Williams, Ronald J. (1992). “Simple statistical gradient-following algorithms for connectionist reinforcement learning”. In: *Mach. Learn.* 8.3-4, pp. 229–256. ISSN: 0885-6125. DOI: 10.1007/BF00992696. URL: <http://link.springer.com/10.1007/BF00992696>.