

# Transforming Worlds: Automated Involutive MCMC for Open-Universe Probabilistic Models

George Matheos<sup>1,2,\*</sup>

Alexander K. Lew<sup>1,\*</sup>

Matin Ghavamizadeh<sup>1,2</sup>

Stuart Russell<sup>2</sup>

Marco Cusumano-Towner<sup>1</sup>

Vikash K. Mansinghka<sup>1</sup>

<sup>1</sup>MIT <sup>2</sup>UC Berkeley \*Equal contribution

GEORGEMATHEOS@BERKELEY.EDU

ALEXLEW@MIT.EDU

MGHAVAMI@BERKELEY.EDU

RUSSELL@BERKELEY.EDU

MARCOCT@MIT.EDU

VKM@MIT.EDU

## 1. Introduction

Open-universe probabilistic models (OUPMs) enable Bayesian inference about the existence and attributes of latent objects underlying observed data, as well as their interconnections. Prominent applications include seismic monitoring for global nuclear safety (Arora et al., 2013; Arora, 2012); information extraction from natural language documents (Russell et al., 2016); generating realistic 3D scenes with varying numbers of objects (Yeh et al., 2012); “inverse graphics” approaches to breaking CAPTCHAs (Mansinghka et al., 2013) and inferring 3D scenes from 2D data (Kulkarni et al., 2015; Zinberg et al., 2019); and simultaneous de-duplication, cleaning, and record linkage from real-world databases with millions of records (Lew et al., 2020). Because the state space in OUPMs has a priori unknown dimension, popular black-box inference algorithms such as Hamiltonian Monte Carlo cannot be applied. Instead, applications rely on custom MCMC kernels. These kernels use application-specific, data-driven heuristics to intelligently delete, split, or merge hypothesized objects, “birth” new objects, and modify the properties of and relationships between objects. Although these proposals can be designed and justified in the reversible-jump MCMC framework, this is often challenging. As Brooks et al. (2003) noted, “the application of reversible jump...has predominantly remained within the domain of the MCMC expert,” due to the difficulty of deriving and implementing effective RJMCMC kernels.

This abstract introduces a new framework for OUPM inference that makes it easier to design and implement custom, data-driven kernels. We adapt involutive MCMC (Cusumano-Towner et al., 2020; Neklyudov et al., 2020), a generalization of reversible-jump MCMC, to the setting of OUPMs. This enables software to automatically generate correct and scalable implementations of complex, application-specific kernels from high-level user specifications. Users design inference programs that propose incremental changes to possible worlds, creating, deleting, or modifying objects according to data-driven, application-specific heuristics; these proposals are automatically converted into stationary MCMC kernels via an accept/reject step. To automatically compute the acceptance probability, our approach leverages program tracing and dependency tracking (for efficient computation of proposal and model densities) and differentiable programming (for the Jacobian determinant). Preliminary experiments on (1) an auditory scene analysis application (Cusumano et al., 2018) (Figure 3) and (2) Gaussian mixture modeling with an unknown number of components (Richardson and Green, 1997) (Figure 4) show that our approach enables data-driven kernels that

are faster than generic probabilistic programming language (PPL) inference and generic birth/death RJMCMC kernels without application-specific customizations.

**Related Work.** We build on involutive MCMC, introduced by Cusumano-Towner et al. (2020) as a framework supporting higher degrees of automation for application-specific MCMC algorithms, and (independently) as a mathematical generalization of many classic MCMC methods by Neklyudov et al. (2020). Our main theoretical contribution is a generalization of involutive MCMC to OUPMs. We compare our approach to the generic ancestral resampling MH algorithm (Wingate et al., 2011) used by many probabilistic programming languages supporting OUPMs, such as BLOG (Milch et al., 2005a). Previous work has explored using application-specific kernels for BLOG models (Milch and Russell, 2012), but these require modifying the BLOG inference engine’s source code and manually deriving parts of the acceptance ratio. The Gen (Cusumano-Towner et al., 2019) and Stochaskell (Roberts et al., 2019) PPLs support forms of automated reversible-jump MCMC, but neither features high-level OUPM inference constructs; custom kernels in both languages operate on low-level program traces, not high-level object-based world representations, which can compromise both ease of use and asymptotic performance. Furthermore, the automation technique proposed by Roberts et al. (2019) does not support unrestricted use of auxiliary variables, and (e.g.) cannot handle the algorithms benchmarked in Figure 3.

## 2. Open Universe Probabilistic Models

Figure 1 shows an open-universe model for inferring the set of seismic events underlying observed detections from seismic monitoring stations. A priori, it is unknown how many events there are, and which detections correspond to the same event (or, in the case of spurious false positives, to no event at all). Standard probabilistic graphical models are ill-suited for encoding this sort of uncertainty about the *number* of objects reflected in a dataset and their causal relationships. *Open-universe probabilistic models* enable reasoning about such problems, by defining generative processes over entire relational domains.

**Model specification.** An OUPM is fully specified by: (1) a finite set  $\mathcal{T}$  of object types (in Figure 1, `Event`, `Station`, and `Detection`); (2) a finite set  $\mathcal{N}_\tau$  of *possible origins* for each object type  $\tau$  (in Figure 1, a `Detection` may originate from a `(Station, Event)` pair, or from a single `(Station)` that records a false positive); (3) a countable set  $\mathcal{P}$  of typed *object properties*  $P(\tau_1, \dots, \tau_n)$  (such as `magnitude(Event)` and `reading(Detection)`); and (4) a *contingent Bayesian network*  $\mathcal{C}$  over a set of random variables that encode the number of objects that exist, as well as their properties and relationships (Milch et al., 2005b).

**Contingent Bayesian networks.** The probabilistic structure of an OUPM is given by a *contingent Bayesian network*  $\mathcal{C}$  over an infinite set  $\mathcal{V}$  of *possible variables*. These variables are defined in terms of *objects*: we write  $\tau((o_1, \dots, o_n), i)$  to refer to the  $i^{\text{th}}$  object of type  $\tau$  with origin  $(o_1, \dots, o_n)$ . Then each node  $v \in \mathcal{V}$  of the contingent Bayesian network corresponds to either a *number variable*  $N_\tau(o_1, \dots, o_n)$ , which represents the number of objects of type  $\tau$  with origin  $(o_1, \dots, o_n)$ , or to a *property variable*  $P(o_1, \dots, o_n)$ , which represents the value of the property  $P$  for the given tuple of objects. Like an ordinary Bayesian network, a contingent Bayesian network uses a directed graph to describe the probabilistic dependencies among variables. However, the edges in a CBN’s graph are contingent: they are

```

1  type Event, Station, Detection
2  property function amplitude(e :: Event)
3      return amp ~ exponential(100.0)
4  end
5  function is_false_positive(d::Detection)
6      return length(origin(d)) == 1
7  end
8  property function magnitude(d :: Detection)
9      if is_false_positive(d)
10         reading ~ normal(1, 0.5)
11     else
12         (station, event) = origin(d)
13         richter_mag = log_10(get(amplitude(event)))
14         reading ~ normal(richter_mag, 0.2)
15     end
16     return reading
17 end

18 number function Station()
19     return 2
20 end
21 number function Event()
22     return num ~ poisson(5)
23 end
24 number function Detection(:: Station)
25     return num ~ poisson(4)
26 end
27 number function Detection(:: Station, :: Event)
28     return num ~ Int(bernoulli(0.8))
29 end

30 observation_model function detections()
31     detections = get_object_set(Detection)
32     return [get(magnitude(d)) for d in detections]
33 end
    
```

Figure 1: An OUPM for seismic monitoring, inspired by Arora et al. (2013), in which **Detections** are received at seismic monitoring **Stations**. Each **Station** receives some number of false-positive **Detections**, and for each **Event**, each **Station** has probability 0.8 of detecting it. Each **Event** has an **amplitude**, and a reading of a Richter-scale **magnitude** is recorded for each **Detection**, which is distributed around the underlying **Event**’s magnitude when the **Detection** is not a false-positive.

---

### Algorithm 1 Automated Involutive MCMC for OUPMs

---

```

procedure OUPM-IMCMC( $p, q, f, x$ )
    ( $y, q_x(y)$ )  $\leftarrow$  SAMPLEANDSCORE( $q, x$ )
    ( $U, M, y', J$ )  $\leftarrow$  PROCESSINVOLUTION( $f, (x, y)$ )
    ( $x', \frac{p(x')}{p(x)}$ )  $\leftarrow$  UPDATEWORLD( $x, U, M$ )
     $q_{x'}(y') \leftarrow$  SCORE( $q, x', y'$ )
    with probability  $\frac{p(x')}{p(x)} \times q_{x'}(y') \times \frac{1}{q_x(y)} \times J$  return  $x'$  else return  $x$ 
end procedure

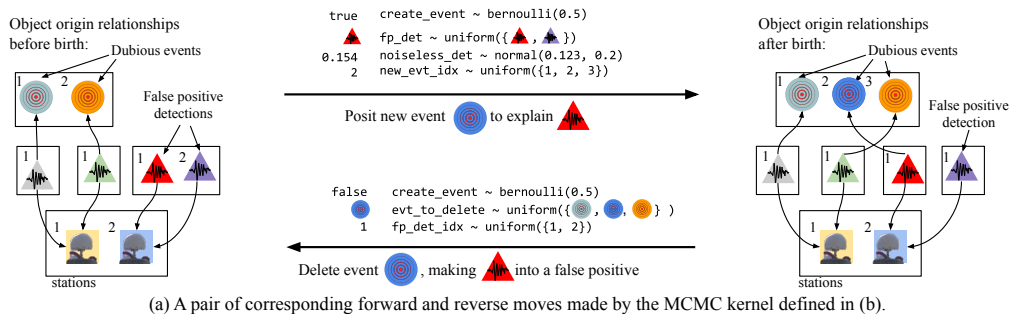
procedure PROCESSINVOLUTION( $f, (x, y)$ )
     $\triangleright$  Track continuous reads and writes, manipulation moves  $M$ , and property updates  $U$ 
    ( $\text{Rd}, \text{Wr}, M, U, y'$ )  $\leftarrow$   $\{\}, \{\}, [], \{\}, \{\}$ 
    Execute  $f$ , but with
        each manipulation command  $m$  (create, delete, change, split, merge)  $\equiv (M \leftarrow M \cup \{m\})$ 
        "proposed[ $k$ ]"  $\equiv$  (if  $y[k]$  is continuous:  $\text{Rd} \leftarrow \text{Rd} \cup \{y[k]\}; y[k]$ )
        "set backward[ $k$ ] =  $v$ "  $\equiv$  (if  $v$  is continuous:  $\text{Wr} \leftarrow \text{Wr} \cup \{y[k]\}; y' \leftarrow y' \cup \{k \mapsto v\}$ )
        "get( $P(o_1, \dots, o_n)$ )"  $\equiv (v \leftarrow x[P(o_1, \dots, o_n)];$  if  $v$  is continuous:  $\text{Rd} \leftarrow \text{Rd} \cup \{v\}; v$ )
        "set  $P(o_1, \dots, o_n) = v$ "  $\equiv$  (if  $v$  is continuous:  $\text{Wr} \leftarrow \text{Wr} \cup \{v\}; U[P(o_1, \dots, o_n)] \leftarrow v$ )
     $J \leftarrow$  uninitialized  $|\text{Rd}| \times |\text{Rd}|$  matrix
    for  $i$  in  $\{1, \dots, |\text{Wr}|\}$ :
         $\triangleright$  Execute  $f$  with reverse-mode AD to compute  $\frac{\partial \text{Wr}[i]}{\partial \text{Rd}[j]}$  for each  $j \in \{1, \dots, |\text{Rd}|\}$ 
         $J[:, i] \leftarrow \nabla_{\text{Rd}}(\text{Wr}[i])$ 
    return ( $U, M, y', |\det J|$ )
end procedure
    
```

---

labeled with predicates, involving the other possible variables, determining the conditions under which they are *active*. The conditional distribution for a variable  $v$ ,  $p_v(\omega[v] \mid \omega_{Pa_\omega(v)})$ , may depend only on  $v$ ’s *active* parents  $Pa_\omega(v)$  in a particular possible world  $\omega$ .

### Induced probability distribution over minimal self-supporting instantiations.

Fixing a set  $U \subseteq \mathcal{V}$  of *observed* variables (in Figure 1,  $\{\text{detections}()\}$ ), the CBN  $\mathcal{C}$  can be used to define a probability distribution over *minimal self-supporting instantiations* (MSSIs) for  $U$ : *finite* assignments  $w$  to a *subset* of variables  $\text{vars}(w) \subseteq \mathcal{V}$ , where  $U \subseteq \text{vars}(w)$ , such that (1)  $w$  is *self-supporting*—if  $v \in \text{vars}(w)$  and  $u \in Pa_w(v)$  then  $u \in \text{vars}(w)$ ; and (2)  $w$  is *minimal*—if any variable  $v \notin U$  were removed from  $w$  it would no longer be self-supporting.



(a) A pair of corresponding forward and reverse moves made by the MCMC kernel defined in (b).

```
is_detection_of(e, d) = !is_false_positive(d) && origin(d)[2] == e
is_dubious(e, ds) = sum(is_detection_of(e, d) for d in ds) == 1
```

```
gen function birth_death_proposal(prev_world)
  dets, events, stations = [get_object_set(prev_world, T)
    for T in [Detection, Event, Station]]

  # decide whether to create or delete an event
  false_positives = filter(is_false_positive, dets)
  dubious_events = filter(e -> is_dubious_event(e, dets), events)
  create_prob = isempty(false_positives) ? 0.0 :
    isempty(dubious_events) ? 1.0 : 0.5
  create_event ~ bernoulli(create_prob)

  # make additional choices based on move type
  if create_event
    fp_det ~ uniform(false_positives)
    # what would the detection of the new event be with no noise?:
    noiseless_det ~ normal(get(prev_world, :magnitude, fp_det), 0.2)
    new_evt_idx ~ uniform(range(1, length(events) + 1))
  else
    evt_to_delete ~ uniform(dubious_events)
    detection = find(d -> is_detection_of(evt_to_delete, d), dets)
    num_fps = sum(first(origin(d)) == first(origin(detection))
      for d in fp_detections)
    fp_det_idx ~ uniform(range(1, num_fps + 1))
  end
end
```

```
involution function birth_death_involution()
  set backward[:create_event] = !proposed[:create_event]

  if proposed[:create_event]
    # create the new dubious event
    new_event = create Event with index proposed[:new_evt_idx]
    new_amp = 10^(proposed[:noiseless_det]) # Richter -> amplitude
    set amplitude(new_event)[:amp] = new_amp

    # change the false positive to be a true detection of the new event
    station = origin(proposed[:fp_det])
    change proposed[:fp_det] to origin(station, new_event) index 1

    # fill in details of the backward move
    set backward[:evt_to_delete] = new_event
    set backward[:fp_detection_idx] = index(proposed[:fp_detection])
  else
    # delete the dubious event
    dubious_evt = proposed[:evt_to_delete]
    delete dubious_evt

    # find detection of this event and change to a false positive
    dets = get_object_set(Detection)
    det = find(d -> is_detection_of_event(dubious_evt, d), dets)
    station = (origin(det))[1]
    new_fp = change det to origin(station, ) index proposed[:fp_det_idx]

    # fill in details of backward move
    set backward[:fp_det] = new_fp
    set backward[:noiseless_det] = log_10(get(amplitude(dubious_evt)))
    set backward[:new_evt_idx] = index(proposed[:evt_to_delete])
  end
end
```

(b) An involutive MCMC kernel consisting of a proposal generative function and an involution.

Figure 2: (a) An execution of a MCMC kernel transforming the state of the seismic monitoring model on the left to that on the right. The move, with random choices shown above the right arrow, posits a new event to explain a false positive detection. The involution sets the new event’s amplitude to  $10^r$ , converting Richter-scale magnitude  $r$  to an event amplitude, so the acceptance probability’s  $(\frac{d(\mu \circ f^{-1})}{d\mu}(x, y))$  term is  $\log(10)10^r$ ; our system calculates this via automatic differentiation. The kernel can also make the reverse move ( $\leftarrow$ ). (b) The involutive MCMC kernel visualized in (a).

Let  $\mathcal{W}$  be the set of all possible minimal self-supporting instantiations for a fixed set  $U$  of observed variables, and let  $\text{vars}(\mathcal{W}) = \{\text{vars}(w) \mid w \in \mathcal{W}\}$  be the set of possible sets-of-variables that can be jointly specified by a minimal self-supporting instantiation. For  $\mathbf{v} \in \text{vars}(\mathcal{W})$ , let  $\mathbf{v}_d$  be the set of all possible joint assignments to the *discrete* variables in  $\mathbf{v}$ . Then we can define a reference measure  $\mu$  over the space  $\mathcal{W}$  of MSSIs, given by

$$\mu(E) = \sum_{\mathbf{v} \in \text{vars}(\mathcal{W})} \sum_{w_d \in \mathbf{v}_d} \lambda^{|\mathbf{v}| - |w_d|} (\{w'_c \mid w' \in E \wedge \text{vars}(w') = \mathbf{v} \wedge \forall v \in \text{vars}(w_d), w'[v] = w_d[v]\}),$$

where  $w_c$  is the vector of values of the continuous variables in an MSSI  $w$ , and  $\lambda^k$  is the Lebesgue measure on  $\mathbb{R}^k$ . With respect to  $\mu$ , the probability density for the OUPM, as a distribution over MSSIs, is given by  $p(w) = \prod_{v \in \text{vars}(w)} p_v(w[v] \mid \{u \mapsto w[u] \mid u \in Pa_w(v)\})$ .<sup>1</sup>

1. Milch et al. (2005b)’s formulation of contingent Bayesian networks does not support continuous random variables. See Appendix B for a more careful development of CBNs, and an extension (to our knowledge, novel) to the continuous-variable case.



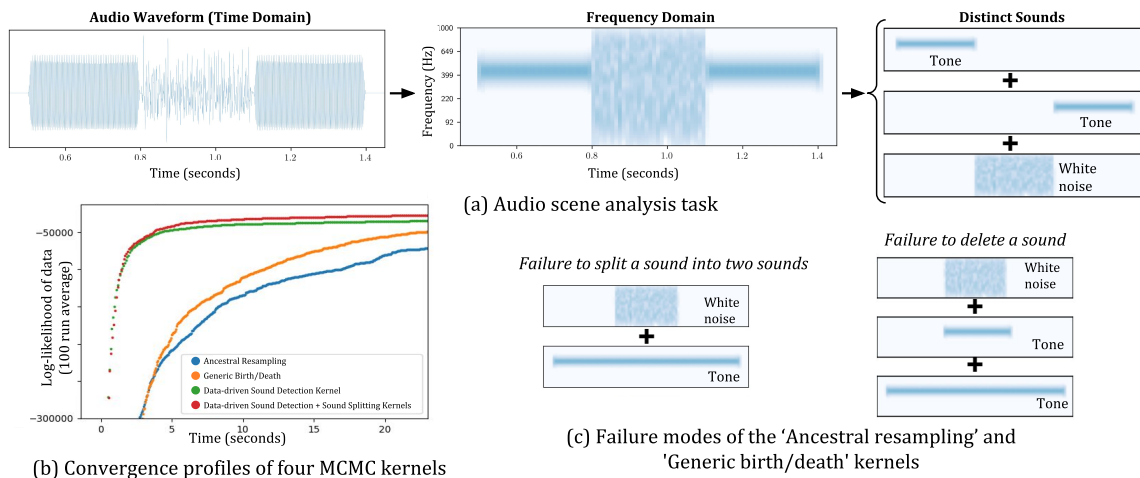


Figure 3: An OUPM for auditory perception. (a) A waveform (left) is observed as a spectrogram in the frequency domain (center), from which we infer the number and type of constituent sound sources that together explain the observed data (right). (b) Convergence profiles of 4 MCMC kernels; the two data-driven iMCMC kernels outperform generic PPL and generic RJMCMC kernels. (c) Approximate posterior samples reflecting typical failure modes of non-data-driven kernels.

### 3. Involutive MCMC for Open Universe Models

An involutive MCMC kernel (Cusumano-Towner et al., 2020; Neklyudov et al., 2020) for a target density  $p$  on a measure space  $(X, \Sigma_X, \mu_X)$  uses a three-step process to transform a state  $x$  into a next state  $x'$ . First, auxiliary variables  $y$  are sampled from a state-dependent proposal with density  $q_x$ , defined over a measure space  $(Y, \Sigma_Y, \mu_Y)$ . Second, an *involution*<sup>2</sup>  $f$  is applied to  $(x, y)$  to obtain  $(x', y')$ :  $x'$  is the proposed next state, and  $y'$  is an auxiliary variable value that would cause the kernel to propose the old state  $x$  from the new state  $x'$  (since  $f(x', y') = (x, y)$ ). Third, an acceptance probability  $\alpha(x, y, x', y')$  is computed. With probability  $\alpha$ ,  $x'$  is the Markov chain's next state; otherwise the state  $x$  is repeated. Cusumano-Towner et al. (2020) show that this process yields a stationary kernel for  $p$  when

$$\alpha(x, y, x', y') = \frac{p(x')q_{x'}(y')}{p(x)q_x(y)} \times \left( \frac{d(\mu \circ f^{-1})}{d\mu}(x, y) \right), \quad (1)$$

where  $\frac{d(\mu \circ f^{-1})}{d\mu}$  is the Radon-Nikodym derivative of the pushforward of  $\mu = \mu_X \times \mu_Y$  through  $f$  with respect to  $\mu$ . In general, this does not yield a constructive procedure for *computing* this Radon-Nikodym derivative, but Cusumano-Towner et al. (2020) show that when  $p$  and  $q$  are distributions over a space  $\mathcal{D}$  of finite *dictionaries*, i.e., finite sets of key-value pairs, the Radon-Nikodym derivative can be computed as the determinant of a Jacobian matrix. It turns out (Theorem 1) that a similar result holds when  $p$  is an open-universe model. Suppose  $f$  is an involution on the space  $\mathcal{W} \times \mathcal{D}$ , and that there is a countable partition  $\{C_i\}$  of  $\mathcal{W} \times \mathcal{D}$  such that: (i)  $(x, y) \in C_i \implies f(x, y) \in C_i$ , (ii) for all  $(x, y) \in C_i$ , the total number of continuous values in minimal self-supporting instantiation  $x$ , plus the number of continuous values in dictionary  $y$ , is constant, and (iii) for each  $C_i$ , for  $(x, y) \in C_i$ , there

2. An involution is a function that is its own inverse:  $f(f(x, y)) = (x, y)$ .

is a differentiable bijection  $h_i : \mathbb{R}^{k_i} \rightarrow \mathbb{R}^{k_i}$  mapping the continuous values in  $x$  and  $y$  to continuous values in  $x'$  and  $y'$  (where  $(x', y') = f(x, y)$ ). Then, we have:

**Theorem 1** *For any open universe model with density  $p$ , any collection of distributions over finite dictionaries with densities  $q_x$ , and any involution satisfying the above conditions, involutive MCMC with  $\alpha(x, y, x', y') = \frac{p(x')q_{x'}(y')}{p(x)q_x(y)} \times |\det(Jh_i(x, y))|$  is stationary for  $p$ , where  $Jh_i(x, y)$  is the Jacobian of  $h_i$  at the continuous variables of  $(x, y)$ , and  $i$  is s.t.  $(x, y) \in C_i$ .*

The proof involves a translation to the finite dictionary case, and exploits the structure of minimal self-supporting instantiations to guarantee the necessary properties of  $f$ .

#### 4. Automation of Acceptance Probability via PPLs and AD

We have extended the Gen probabilistic programming system (Cusumano-Towner et al., 2019) with new constructs for OUPMs (sample code in Figure 1) and automated involutive MCMC for these models (Figure 2, Algorithm 1). The contributions are: (1) a new modeling language for open-universe models, capable of efficiently simulating MSSIs, evaluating densities, and modifying world states; (2) a new DSL for defining object-based involutive MCMC proposals (Appendix C); and (3) an *automated* implementation of involutive MCMC for open-universe models, which automatically calculates acceptance ratios via probabilistic and differentiable programming techniques (Algorithm 1). See Appendix D for more details.

#### 5. Experiments

**Auditory scene analysis.** We implemented a variant of Cusumano et al. (2018)’s model of human auditory scene perception, which is based on an OUPM of sound sources with properties including amplitude, type, pitch, and duration. See Figure 3 for an example problem, and Appendix A.2 for details. We implemented two custom, data-driven kernels: a **data-driven sound detection kernel** (Figure 11) that applies image-filtering and edge-detection to frequency-domain spectrograms and proposes new audio sources to explain sounds it detects in the data but not in the currently-hypothesized world, and a **data-driven sound splitting kernel** (Figure 12) that splits long sounds into shorter ones at endpoints detected in the observed soundwaves. We also implemented two baseline inference strategies: an **ancestral resampling** kernel (Figure 9) based on generic PPL inference in BLOG; and a **generic birth/death** RJMCMC kernel (Figure 10) that proposes to create a sound with parameters generated from the prior, or delete a sound. Figure 3(b) shows quantitative performance comparisons indicating superior performance for custom data-driven kernels. Common failure modes of the generic kernels are illustrated in Figure 3(c).

**Open-universe Gaussian mixture modeling.** We also implemented Richardson and Green (1997)’s data-driven inference algorithm for clustering in a Gaussian mixture model with an unknown number of components. See Appendix A.1 for details, including an illustration in Figure 4. We have found that BLOG inference typically fails to accept moves to world states with different numbers of clusters, whereas the data-driven algorithm rapidly mixes. Note that Richardson and Green (1997)’s algorithm includes a split-merge kernel with a Jacobian term that our approach calculates automatically.

## References

- Nimar Arora, Stuart Russell, and Erik Sudderth. Net-visa: Network processing vertically integrated seismic analysis. *The Bulletin of the Seismological Society of America*, 103: 709–729, 04 2013. doi: 10.1785/0120120107.
- Nimar S Arora. *Model-based Bayesian Seismic Monitoring*. PhD thesis, EECS Department, University of California, Berkeley, May 2012. URL <http://www2.eecs.berkeley.edu/Pubs/TechRpts/2012/EECS-2012-125.html>.
- Albert S Bregman. *Auditory scene analysis: The perceptual organization of sound*. MIT press, 1994.
- Stephen Brooks, Paolo Giudici, and G.O. Roberts. Efficient construction of reversible jump proposal distributions (with discussion). *J. R. Statist. Soc. B*, 65:3–55, 01 2003.
- Maddie Cusimano, Luke B Hewitt, Josh Tenenbaum, and Josh H McDermott. Auditory scene analysis as bayesian inference in sound source models. In *CogSci*, 2018.
- Marco Cusumano-Towner, Alexander K. Lew, and Vikash K. Mansinghka. Automating involutive mcmc using probabilistic and differentiable programming. 2020.
- Marco F Cusumano-Towner, Feras A Saad, Alexander K Lew, and Vikash K Mansinghka. Gen: a general-purpose probabilistic programming system with programmable inference. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 221–236. ACM, 2019.
- T. D. Kulkarni, P. Kohli, J. B. Tenenbaum, and V. Mansinghka. Picture: A probabilistic programming language for scene perception. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4390–4399, 2015. doi: 10.1109/CVPR.2015.7299068.
- Alexander K. Lew, Monica Agrawal, David Sontag, and Vikash K. Mansinghka. Pclean: Bayesian data cleaning at scale with domain-specific probabilistic programming, 2020.
- Vikash K Mansinghka, Tejas D Kulkarni, Yura N Perov, and Josh Tenenbaum. Approximate bayesian image interpretation using generative probabilistic graphics programs. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 1520–1528. Curran Associates, Inc., 2013. URL <http://papers.nips.cc/paper/4881-approximate-bayesian-image-interpretation-using-generative-probabilistic-graphics-pdf>.
- Brian Milch and Stuart J. Russell. General-purpose MCMC inference over relational structures. *CoRR*, abs/1206.6849, 2012. URL <http://arxiv.org/abs/1206.6849>.
- Brian Milch, Bhaskara Marthi, Stuart Russell, David Sontag, Daniel L. Ong, and Andrey Kolobov. BLOG: Probabilistic models with unknown objects. In *Proceedings of the Nineteenth International Joint Conference on Artificial Intelligence, IJCAI 2005*, pages 1352–1359. Morgan Kaufmann Publishers Inc., 2005a.

- Brian Milch, Bhaskara Marthi, David Sontag, Stuart Russell, Daniel L. Ong, and Andrey Kolobov. Approximate inference for infinite contingent bayesian networks. In *(AISTATS'05) 10th International Workshop on Artificial Intelligence and Statistics*, January 2005b. URL <https://www.microsoft.com/en-us/research/publication/approximate-inference-infinite-contingent-bayesian-networks/>.
- Brian Christopher Milch. *Probabilistic Models with Unknown Objects*. PhD thesis, EECS Department, University of California, Berkeley, Dec 2006. URL <http://www2.eecs.berkeley.edu/Pubs/TechRpts/2006/EECS-2006-174.html>.
- Kirill Neklyudov, Max Welling, Evgenii Egorov, and Dmitry Vetrov. Involutive MCMC: A Unifying Framework. *arXiv preprint arXiv:2006.16653*, 2020.
- Sylvia Richardson and Peter J. Green. On bayesian analysis of mixtures with an unknown number of components (with discussion). *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 59(4):731–792, 1997. doi: 10.1111/1467-9868.00095. URL <https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/1467-9868.00095>.
- David A. Roberts, Marcus Gallagher, and Thomas Taimre. Reversible jump probabilistic programming. volume 89 of *Proceedings of Machine Learning Research*, pages 634–643. PMLR, 16–18 Apr 2019. URL <http://proceedings.mlr.press/v89/roberts19a.html>.
- S. Russell, O. Lassen, J. Uang, and W. Wang. The physics of text: Ontological realism in information extraction. In *AKBC@NAACL-HLT*, 2016.
- Wei Wang and Stuart Russell. A smart-dumb/dumb-smart algorithm for efficient split-merge mcmc. In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*, UAI'15, page 902–911, Arlington, Virginia, USA, 2015. AUAI Press. ISBN 9780996643108.
- David Wingate, Andreas Stuhlmüller, and Noah Goodman. Lightweight implementations of probabilistic programming languages via transformational compilation. In Geoffrey Gordon, David Dunson, and Miroslav Dudík, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 770–778, Fort Lauderdale, FL, USA, 11–13 Apr 2011. JMLR Workshop and Conference Proceedings. URL <http://proceedings.mlr.press/v15/wingate11a.html>.
- Yi Wu, Siddharth Srivastava, Nicholas Hay, Simon Du, and Stuart J. Russell. Discrete-continuous mixtures in probabilistic programming: Generalized semantics and inference algorithms. *CoRR*, abs/1806.02027, 2018. URL <http://arxiv.org/abs/1806.02027>.
- Yi-Ting Yeh, Lingfeng Yang, Matthew Watson, Noah D. Goodman, and Pat Hanrahan. Synthesizing open worlds with constraints using locally annealed reversible jump mcmc. *ACM Trans. Graph.*, 31(4), July 2012. ISSN 0730-0301. doi: 10.1145/2185520.2185552. URL <https://doi.org/10.1145/2185520.2185552>.
- Ben Zinberg, Marco Cusumano-Towner, and Vikash K. Mansinghka. Structured differentiable models of 3D scenes via generative scene graphs. In *Workshop on Perception*

*as Generative Reasoning*, 2019. URL <https://pgr-workshop.github.io/img/PGR025.pdf>.

## Appendix A. Experimental Details

### A.1. Open-universe Gaussian Mixture Model.

Richardson and Green (1997) develop a series of highly data-driven MCMC kernels for a Gaussian mixture model with an unknown number of mixture components. Here we compare their algorithm, implemented using our system for involutive MCMC in OUPMs, to the generic inference engine in BLOG.

The model is as follows. The number of mixture components,  $k$ , is a Poisson random variable. The mixture weights and the parameters of each mixture component are drawn from the appropriate conjugate priors:

$$w_1, \dots, w_k \sim \text{Dirichlet}(\delta \mathbf{1}_k) \quad \mu_j \sim \mathcal{N}(\xi, 1/\kappa) \quad \sigma_j^2 \sim \text{inverse-gamma}(\alpha, \beta),$$

where  $\mathbf{1}_k$  denotes a vector of ones in  $\mathbb{R}^k$ ;  $w_j$ ,  $\mu_j$ , and  $\sigma_j^2$  denote the weight, mean, and variance of the  $j$ -th mixture component respectively; and  $\alpha$ ,  $\beta$ ,  $\delta$ ,  $\kappa$ , and  $\xi$  are fixed hyper-parameters. Lastly, for the  $i$ -th data point we sample a cluster allocation, and given the cluster allocation, we sample a value from the assigned cluster; ie.

$$z_i \sim \text{Categorical}(w_1, \dots, w_k) \quad y_i \sim \mathcal{N}(\mu_{z_i}, \sigma_{z_i}^2)$$

Our involutive MCMC proposal consists of first updating the weights, means, variances, and allocations using Gibbs updates, and then performing a split or merge move to change the number of clusters. The Gibbs updates are given by

$$\begin{aligned} w_1, \dots, w_k | \dots &\sim \text{Dirichlet}(w_1 + n_1, \dots, w_k + n_k) \\ \mu_j | \dots &\sim \mathcal{N}\left(\frac{\kappa \xi + \sigma_j^{-2} \sum_{z_i=j} y_i}{\sigma_j^{-2} n_j + \kappa}, (\sigma_j^{-2} n_j + \kappa)^{-1}\right) \\ \sigma_j^2 | \dots &\sim \text{inverse-gamma}\left(\alpha + \frac{1}{2} n_j, \beta + \sum_{z_i=j} (y_i - \mu_i)^2\right) \\ p(z_i = j | \dots) &\propto \mathcal{N}(\mu_j, \sigma_j^2), \end{aligned}$$

where  $n_j$  denotes the number of data points allocated to the  $j$ -th cluster. Our merge moves consist of choosing indices  $j_1 \neq j_2$  and  $j_*$  uniformly at random and merging the  $j_1$ -th and  $j_2$ -th clusters into the  $j_*$ -th cluster. We ensure that the zeroth, first, and second moments of the newly generated cluster match the corresponding moments in the mixture of the original clusters; i.e.

$$\begin{aligned} w_* &= w_1 + w_2 \\ w_* \mu_* &= w_1 \mu_1 + w_2 \mu_2 \\ w_* (\mu_* + \sigma_*^2) &= w_1 (\mu_1 + \sigma_1^2) + w_2 (\mu_2 + \sigma_2^2). \end{aligned} \tag{2}$$

To ensure ergodicity, our split procedure relies on three auxiliary variables

$$u_1 \sim \text{beta}(2, 2) \quad u_2 \sim \text{beta}(2, 2) \quad u_3 \sim \text{beta}(1, 1),$$



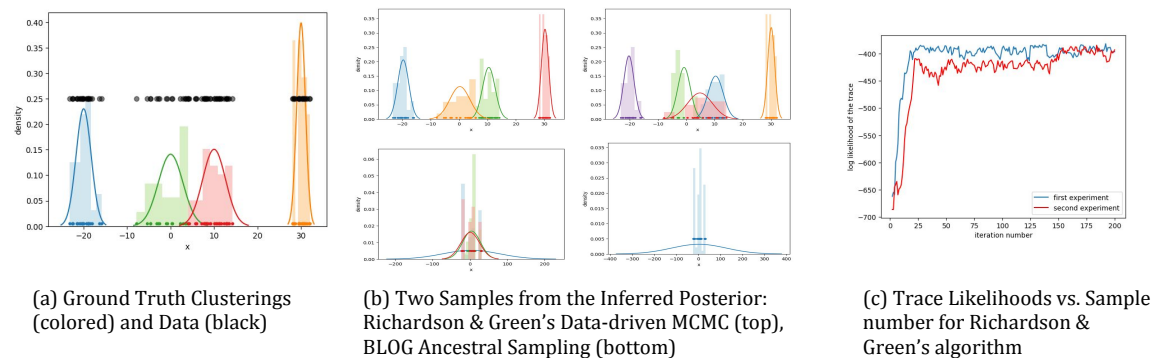


Figure 4: Gaussian Mixture Model Inference Experiment Results

and generates new clusters according to

$$\begin{aligned}
 w_1 &= w_* u_1 & w_2 &= w_*(1 - u_1) \\
 \mu_1 &= \mu_* - u_2 \sigma_* \sqrt{w_2/w_1} & \mu_2 &= \mu_* + u_2 \sigma_* \sqrt{w_1/w_2} \\
 \sigma_1^2 &= u_3(1 - u_2^2) \sigma_*^2 w_*/w_1 & \sigma_2^2 &= (1 - u_3)(1 - u_2^2) \sigma_*^2 w_*/w_2.
 \end{aligned} \tag{3}$$

Using (2) and (3), we can derive our involution's Jacobian matrix (leaving out rows and columns which merely copy a value):

$$J = \begin{matrix} & \begin{matrix} \mu_* & \sigma_*^2 & w_* & u_1 & u_2 & u_3 \end{matrix} \\ \begin{matrix} w_1 \\ \mu_1 \\ \sigma_1^2 \\ w_2 \\ \mu_2 \\ \sigma_2^2 \end{matrix} & \begin{bmatrix} 0 & 0 & u_1 & w_* & 0 & 0 \\ 1 & -\frac{u_2 \sqrt{w_2/w_1}}{2\sigma_*} & 0 & -\frac{u_2 \sigma_* \sqrt{-u_1-1}}{2u_1(u_1-1)} & -\sigma_* \sqrt{w_2/w_1} & 0 \\ 0 & -\frac{u_3 w_* (u_2^2-1)}{w_1} & 0 & \frac{u_3 \sigma_*^2 (u_2^2-1)}{u_1^2} & -\frac{2u_2 u_3 w_* \sigma_*^2}{w_1} & -\frac{w_* \sigma_*^2 (u_2^2-1)}{w_1} \\ 0 & 0 & 1 - u_1 & -w_* & 0 & 0 \\ 1 & \frac{u_2 \sqrt{w_1/w_2}}{2\sigma_*} & 0 & -\frac{u_2 \sigma_* \sqrt{-u_1-1}}{2u_1(u_1-1)} & \sigma_* \sqrt{w_1/w_2} & 0 \\ 0 & \frac{w_* (u_2^2 u_3 - u_2^2 - u_3 + 1)}{w_2} & 0 & \frac{\sigma_*^2 (u_2^2 u_3 - u_2^2 - u_3 + 1)}{(u_1-1)^2} & \frac{2u_2 w_* \sigma_*^2 (u_3-1)}{w_2} & \frac{w_* \sigma_*^2 (u_2^2-1)}{w_2} \end{bmatrix} \end{matrix} .$$

**Given a specification of the inference kernel, our system automatically calculates this Jacobian matrix and its determinant**

$$|\det J| = \frac{w_* |\mu_1 - \mu_2| \sigma_1^2 \sigma_2^2}{u_2(1 - u_2^2) u_3(1 - u_3) \sigma_*^2}$$

without the user needing to work through this derivation.

Figure 4 shows representative results, in which the involutive MCMC proposal accurately recovers the true latent clusters, whereas the ancestral sampling proposal only roughly narrows down the component parameters.

```

type Cluster
type Datapoint
number Cluster() ~ poisson( $\lambda$ ) + 1
number Datapoint() = NUM_DATAPOINTS # global constant
property mean(::Cluster) ~ normal( $\xi, 1/\kappa$ )
property var(::Cluster) ~ inverse_gamma( $\alpha_v, \beta_v$ )
property unnormalized_weight(::Cluster) ~ gamma( $\alpha_w, \beta_w$ )

property function cluster(::Datapoint)
    clust_to_un_weight = Dict{c => get(unnormalized_weight(c))
        for c in get_object_set(Cluster())}
    cluster ~ unnormalized_categorical(clust_to_un_weight)
    return cluster
end

property function value(d::Datapoint)
    cluster = get(cluster(d))
    value ~ normal(get(mean(cluster)), get(var(cluster)))
    return value
end

observation_model function get_datapoints()
    return [get(value(Datapoint with index i))
        for i=1:get_number(Datapoint())]
end
    
```

Figure 5: Gaussian mixture model with an unknown number of components. Instead of imperatively sampling cluster weights  $w_1, \dots, w_n \sim \text{Dirichlet}(\delta \mathbf{1}_k)$ , in our model, each cluster has an *unnormalized cluster weight* property sampled from a Gamma distribution; normalizing these then sampling clusters from the resulting distribution is equivalent to sampling clusters according to weights sampled from a Dirichlet.

```

gen function splitmerge_proposal(prev_world)
    prob_merge = get_number(prev_world, Cluster()) <= 1 ? 0.0 : 0.5
    do_merge ~ bernoulli(prob_merge)
    if is_merge
        {c1, c2} ~ merge_proposal(prev_world)
    else
        {c1, c2} ~ split_proposal(prev_world)
    end
end

gen function merge_proposal(prev_world)
    clusters = get_object_set(prev_world, Cluster())
    merge1 ~ uniform_choice(clusters)
    merge2 ~ uniform_from_list([c for c in clusters if c != merge1])
    new_idx ~ uniform_discrete(1, length(clusters) - 1)
end

gen function split_proposal(prev_world)
    clusters = get_object_set(prev_world, Cluster())
    to_split ~ uniform_choice(clusters)
    idx1 ~ uniform_discrete(1, length(clusters) + 1)
    idx2 ~ uniform_from_list([i for i=1:(length(clusters)+1) if i != idx1])
    u1 ~ beta(2, 2)
    u2 ~ beta(2, 2)
    u3 ~ beta(1, 1)
    mean, var = get(prev_world, :mean, to_split), get(prev_world, :var, to_split)
    mean1, mean2 = mean - u2*sqrt(var * (1-u1)/u1), mean + u2*sqrt(var+u1/(1-u1))
    var1, var2 = u3*(1 - u2^2)*var/u1, (1 - u3)*(1-u2^2)*var/(1-u1)

    for dp in get_object_set(prev_world, Datapoint())
        if get(prev_world, :cluster, dp) == to_split
            y = get(prev_world, :value, dp)
            p = u1 * exp(logpdf(gaussian, y, mean1, var1))
            q = (1-u1) * exp(logpdf(gaussian, y, mean2, var2))
            {to_first => dp} ~ bernoulli(p/(p+q))
        end
    end
end

function merged_component_params(w1, mean1, var1, w2, mean2, var2)
    w = w1 + w2
    mean = (w1*mean1 + w2*mean2) / w
    var = -mean^2 + (w1*(mean1^2 + var1) + w2*(mean2^2 + var2)) / w
    return (w, mean, var)
end

function reverse_split_params(w, mean, var, w1, mean1, var1, w2, mean2, var2)
    u1 = w1/w
    u2 = (mean - mean1) / sqrt(var * w2/w1)
    u3 = var1/var * u1 / (1 - u2^2)
    Return (u1, u2, u3)
end

involution function splitmerge_involution()
    backward[:do_merge] = !proposed[:do_merge]
    if proposed[:do_merge]
        c1, c2 = proposed[:merge1], proposed[:merge2]
        new = merge c1, c2 to index proposed[:new_idx]
        backward[:idx1], backward[:idx2] = index(c1), index(c2)
        backward[:to_split] = new
        mean1, var1, w1 = get(mean(c1)), get(var(c1)), get(unnormalized_weight(c1))
        mean2, var2, w2 = get(mean(c2)), get(var(c2)), get(unnormalized_weight(c2))

        unnorm_weight, mean, var = merged_component_params(mean1, var1, w1,
            mean2, var2, w2)

        set unnormalized_weight(new) = unnorm_weight
        set mean(new) = mean; set var(new) = var

        u1, u2, u3 = reverse_split_params(unnorm_weight, mean, var,
            w1, mean1, var1, w2, mean2, var2)
        backward[:u1], backward[:u2], backward[:u3] = u1, u2, u3

        # reassign datapoints to merged cluster
        for dp in get_object_set(Datapoint())
            if get(cluster(dp)) in [c1, c2]
                set cluster(dp) = new
                backward[:to_first => dp] = (get(cluster(dp)) == c1) ? true : false
            end
        end
    else
        old = backward[:to_split]
        new1, new2 = split old to indices proposed[:idx1], proposed[:idx2]
        backward[:merge1], backward[:merge2] = new1, new2
        backward[:new_idx] = index(old)

        u1, u2, u3 = proposed[:u1], proposed[:u2], proposed[:u3]
        mean, var = get(mean(old)), get(var(old))
        mean1, mean2 = mean - u2*sqrt(var * (1-u1)/u1), mean + u2*sqrt(var+u1/(1-u1))
        var1, var2 = u3*(1 - u2^2)*var/u1, (1 - u3)*(1-u2^2)*var/(1-u1)
        set mean(new1) = mean1
        set var(new1) = var1
        set mean(new2) = mean2
        set var(new2) = var2
        un_weight = get(unnormalized_weight(old))
        set unnormalized_weight(new1) = un_weight * u1
        set unnormalized_weight(new2) = un_weight * (1 - u1)

        # reassign datapoints to split clusters
        for dp in get_object_set(Datapoint())
            if get(cluster(dp)) == old
                set cluster(dp) = backward[:to_first => dp] ? new1 : new2
            end
        end
    end
end
    
```

Figure 6: Data-driven Gaussian mixture model split/merge kernel from [Richardson and Green \(1997\)](#).

## A.2. Auditory Scene Interpretation.

Auditory scene analysis [Bregman \(1994\)](#) is a perception task that involves separating an audio waveform into individual sound sources when the number and properties of the sources are unknown a-priori (Figure 3a). Recently, [Cusimano et al. \(2018\)](#) proposed a computational model of human auditory scene analysis based on an open-universe probabilistic model of sound sources, and showed that probabilistic inference in this model based on MCMC is able to qualitatively reproduce human inferences about auditory scenes, including reporting ambiguity in the same situations as humans.

**The model.** We developed a simplified version of the model of [Cusimano et al. \(2018\)](#) that uses a prior distribution on the number of sounds and the type of each sound (‘tone’ or ‘white noise’). Tones have an unknown frequency, onset time, and duration, and white noise sounds have an unknown amplitude, onset time, and duration.

```

type Source
number Source() ~ uniform_discrete(1, 4)
property is_noise(::Source) ~ bernoulli(0.4)
property freq_in_erb(::Source) ~ uniform(0.4, 24.0)
property amplitude(::Source) ~ normal(8.0, 10.0)
property onset(::Source) ~ uniform(0.1, SCENE_LENGTH)
property duration(::Source) ~ uniform(0.1, SCENE_LENGTH)
property function soundwaves(s::Source)
  ons = get(onset(s)); dur = get(duration(s))
  if get(is_noise(s))
    return render_noise(get(amplitude(s)), ons, dur)
  else
    return render_tone(get(freq(s)), ons, dur)
  end
end
observation_model function observed_sound()
  waves = [get(soundwaves(s)) for s in
            get_object_set(Source)]
  true_sound = sum(waves)
  return observed_sound ~ noisy_matrix(true_sound, 1.0)
end
    
```

Figure 7: Our OUPM for audio inference.

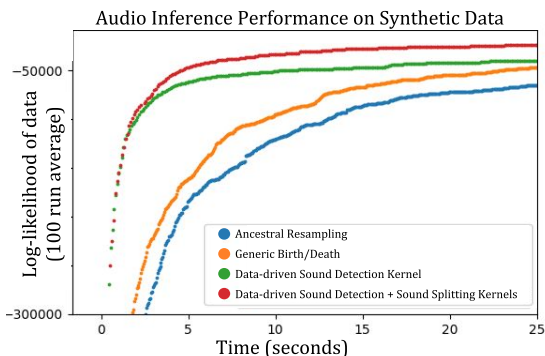


Figure 8: Performance of different inference kernels, measured on 100 randomly-generated audio scenes with 3 underlying sounds.

**Inference Experiments.** We implemented four involutive MCMC kernels for this model in our extended version of Gen. We performed experiments applying these kernels to the audio waveform shown in Figure 3a, which consists of a short tone, followed by a short period of white noise, followed by a short tone at the same frequency as the first tone ([Cusimano et al., 2018](#)). When the noise is very loud, humans tend to hear one continuous tone temporarily obscured by a loud noise, whereas when the noise is quiet, humans hear two distinct short tones, with noise in the middle. Our example uses a softer noise, and our best-performing inference programs, which use application-specific kernels implemented via our framework, recover the true decomposition (two tones separated by noise). However, generic inference kernels often incorrectly infer a single long tone (Figure 3). In addition to testing our kernels on the specific audio scene in Figure 3a, we tested them on 100 audio scenes generated from our model’s prior distribution, with similar results (Figure 8).

**Generic Kernels** The first two kernels we implemented are not application specific or data driven. Our **ancestral resampling** kernel (Figure 9) resimulates a single variable

from the prior distribution (also ancestrally generating any additional variable values which need to be added to the world due to the resimulated variable). Our **generic birth-death** reversible-jump kernel (Figure 10) randomly proposes to create or delete individual sounds, proposing new sounds’ properties from the prior. We cycle this kernel with the ancestral resampling kernel so that property values of existing sounds can more easily be updated. Figure 3(b) shows that in our experiments, these two kernels converge to inferences of latent waveforms under which the observed sounds are substantially less likely than under the latent waveforms our data-driven kernels recover. Figure 3(c) visualizes two common failure mode these kernels exhibit due to their inability to efficiently incorporate evidence from the data, and their inability to make proposals to split long sounds into short ones (which require application-specific logic to calculate endpoints of the new sounds based on the position of the sound being split). We found that the generic birth-death kernel often failed to exit a local minimum in which a single continuous tone is inferred instead of two temporally separated tones of the same frequency (Figure 3c, left). The generic kernel performs worst—it usually infers one continuous tone, but often also infers spurious tones.

**Aside: proposing from the prior in involutive MCMC.** To facilitate proposing values from the prior in involutive MCMC kernels, we have implemented a **regenerate** command which can be called from within the involution DSL. (For examples, see Figures 9, 10, and 11.) The command **regenerate properties of object** samples a value from the prior for each property of the given object (using the prior as it will appear in  $x'$ ). This complicates our picture of involutive MCMC slightly, as there are now random values being sampled during the execution of the involution function, instead of only in the proposal probabilistic program. We thus must view the proposal function as only implementing part of the  $q$  distribution and sampling part of  $y$ , with the rest of  $q$  implemented by the involution DSL’s **regenerate** commands.

To ensure involutiveness, when a transition  $(x, y) \rightarrow (x', y')$  uses a **regenerate** command, the reversing move  $(x', y') \rightarrow (x, y)$  must specify in  $y$  the value sampled by regeneration during the forward move. To specify reverse-move values for non-regenerated values, users use syntax **backward[addr] = value**, where **addr** is the address (variable name) at which the proposal function should sample this value. Since users do not specify an address for the values in  $y$  which are regenerated (unlike values in  $y$  which are sampled in the proposal function), we provide the following additional syntax: **backward[property(object)] = regenerated** sets the value in  $y$  for regenerating the variable **property(object)** to the value of **property(object)** in  $x$ , and **backward[properties of object] = regenerated** does this for all properties of **object**. See Section E for more details.

**Custom kernels.** Some specific domains, like data-clustering, have seen standard MCMC kernels developed for data-driven inference, which are fairly conceptually simple but previously required fairly deriving fairly involved calculations (Richardson and Green, 1997; Wang and Russell, 2015). In addition to facilitating the implementation of these types of kernels and removing the need to manually derive acceptance probabilities, involutive MCMC facilitates the development of data-driven kernels in more complicated domains which involve complicated logic to use information from data effectively. For instance, Cusimano et al. (2018)’s experiments used MCMC inference with a proposal based neural-

```

gen function ancestral_resampling_proposal(prev_world)
  var_to_resample ~ uniform_choice(variables(prev_world))
  if is_number_variable(var_to_resample)
    new_num ~ uniform_discrete(1, 4)
  end
end

involution function ancestral_resampling_involution()
  backward[:var_to_resample] = proposed[:var_to_resample]
  if is_number_variable(proposed[:var_to_resample])
    set number[Source()] = proposed[:new_num]
    backward[:new_num] = get_number(Source())
    for i=(get_number(Source())+1):proposed[:new_num]
      regenerate properties of (Source with index i)
    end
    for i=(proposed[:new_num]+1):get_number(Source())
      backward[properties of (Source with index i)] = regenerated
    end
  else
    regenerate proposed[:var_to_resample]
    backward[proposed[:var_to_resample]] = regenerated
  end
end
    
```

Figure 9: Our **ancestral resampling** kernel, which selects a random variable and samples a new value for it from the prior. The only complication is that when we sample the number variable (changing the number of audio sources in the world), we must also generate new properties for any newly created audio sources. (These properties are also generated from the prior.) For details on generating values from the prior in involutive MCMC, see Section E.

```

gen function generic_birthdeath_proposal(prev_world)
  do_birth ~ bernoulli(0.5)
  sources = get_object_set(prev_world, Source)
  if do_birth
    idx ~ uniform_discrete(1, length(sources) + 1)
  else
    to_delete ~ uniform_choice(sources)
  end
end

involution function generic_birthdeath_involution()
  backward[:do_birth] = !proposed[:do_birth]
  if proposed[:do_birth]
    new = create Source at index proposed[:idx]
    backward[:to_delete] = new
    regenerate properties of new
  end
  old = delete proposed[:to_delete]
  backward[:idx] = index(old)
  backward[properties of old] = regenerated
end
    
```

Figure 10: The **generic birth/death** kernel either creates or deletes a random sound. When a new sound is created, its properties are all generated from the prior. For details on generating values from the prior in involutive MCMC, see Section E.

network transformations of observed soundwaves, implemented in an early version of our inference programming system. For our experiments in this paper, we developed a roughly 250-LOC image-filtering and edge-detection algorithm to analyze observed soundwaves and detect possible errors in the currently inferred latent sounds. We called this sound-analysis program from within 2 involutive MCMC kernels we developed, which our system automated the implementation of, including efficiently updating world states and calculating the Jacobian correction term.

Our **data-driven sound detection kernel** (Figure 11) either proposes to delete a sound, or to create a new sound which with high probability is proposed at a time and frequency where our analysis program detected a sound in the data but not the currently inferred latents. Figure 3(b) and Figure 8 show that this kernel both initially performs better than the generic kernels, and ultimately converges to a better state.

We also add a **data-driven sound splitting kernel** (Figure 12), which either proposes to merge two short tones/noises into a long one, or splits a long sound into two short ones,

```

struct DetectedSound # the perfect `DetectedSound` for the object
  type              # currently at the given index in the world
  amp_or_erb        function det_snd_for_source_at_idx(wrld, idx)
  onset             s = Source with index idx
  duration           is_noise = get(wrld, :is_noise, s),
end
return DetectedSound(is_noise ? :noise : :tone,
struct DeathAction
  type              get(wrld, is_noise ? :amp : :erb, s),
  idx               get(wrld, :onset, s), get(wrld, :duration, s))
end

gen function propose_detected_sound_proposal(prev_world)
  num_sounds = get_number(prev_world, Source)
  # get possible birth and death actions, and heuristic scores for these
  birth_to_score = detect_and_score_possible_births(prev_world)
  death_to_score = score_possible_deaths(prev_world)
  # decide whether to do a birth or death move (create vs delete a sound)
  birthscore = sum(values(birth_to_score))
  deathscore = sum(values(death_to_score))
  birthprior1 = birthscore / (birthscore + deathscore)
  birthprior2 = (birthprior1 + BIRTH_PRIOR * PROB_RANDOMLY_CHOOSE) /
    (1 + PROB_RANDOMLY_CHOOSE)
  birthprior = (num_sounds == 4 || isempty(birth_to_score)) ? 0 : (
    num_sounds == 0 ? 1. : birthprior2)
  do_birth ~ bernoulli(birthprior)
  if do_birth
    # to aid reversibility, there is a chance of a "dumb birth" from prior
    smart_birth_prior = isempty(birth_to_score) ? 0 : SMART_BIRTH_PRIOR
    do_smart_birth ~ bernoulli(smart_birth_prior)
  end

  if !do_birth || do_smart_birth
    # choose type of object to perform create/delete
    action_to_score = do_birth ? birth_to_score : death_to_score
    type_to_score = get_type_to_score(action_to_score)
    objtype ~ unnormalized_categorical(type_to_score)
    typed_action_to_score = Dict{a => s for (a, s) in action_to_score
      if a.type == objtype)
    action ~ unnormalized_categorical(typed_action_to_score)
  end
  if do_birth
    if do_smart_birth
      {:properties} ~ sample_properties(objtype, action)
    end # else: generate object properties from the prior (see involution)
    idx ~ uniform_discrete(1, num_sounds + 1)
  else
    {:rev} ~ sample_death_reversing_randomness(prev_world, objtype, action)
  end
end

gen function sample_properties(objtype, detected::DetectedSound)
  if objtype == noise
    amp ~ normal(detected.amp_or_erb, AMP_STD)
  else
    erb ~ truncated_normal(detected.amp_or_erb, ERB_STD, ERB_RANGE...)
  end
  onset ~ truncated_normal(detected.onset, ERB_STD, ONSET_RANGE...)
  dur ~ truncated_normal(detected.dur, DUR_STD, DUR_RANGE...)
end

gen function sample_death_reversing_randomness(wrld, objtype,
  action::DeathAction)
  # detect sounds in scene but not in latents, ignoring the object being
  # deleted (ie. the detections possible after deleting this sound)
  possible_births = detect_births(wrld, Set(action.idx))
  birth_to_score = Dict{ # score detections by similarity to deleted sound
    b => prob_of_sampling(det_snd_for_source_at_idx(wrld, action.idx),
      b, SCENE_LENGTH) for b in possible_births if b.type == objtype
  }
  reverse_smart_prior = length(birth_to_score) == 0 ? 0. : SMART_BIRTH_PRIOR
  reverse_is_smart ~ bernoulli(reverse_smart_prior)

  if reverse_is_smart
    # which detection is selected to recreate the source we're deleting?
    reverse_birth ~ unnormalized_categorical(birth_to_score)
  end
end
end

const SMART_BIRTH_PRIOR=0.9, PROB_RANDOMLY_CHOOSE_BD=0.1, BIRTH_PRIOR=0.5,
AMP_STD=1.0, ERB_STD=0.5, ONSET_STD=0.1, DUR_STD=0.1, ERB_RANGE=(0.4,
24.), NOISE_PRIOR_PROB=0.4, SCENE_LENGTH = 2., ONSET_RANGE = (0.0,
SCENE_LENGTH), DUR_RANGE = (0.1, 1.)

function detect_and_score_possible_births(prev_world)
  # use edge-detection and image filtering to analyze the observed sounds;
  # return sounds in the observations but not in the latents,
  # ignoring latents in the given set (so in this case, ignoring none)
  detected_births = detect_sounds(prev_world, Set())
  # assign a score to each birth action using a heuristic
  return Dict{b -> birth_score(b, prev_world) for b in detected_births}
end

function score_possible_deaths(prev_world)
  return Dict{DeathAction(
    get(prev_world, :is_noise, source) ? :noise : :tone, index(source)
  ) => death_score(source, prev_world) # heuristic scoring function
    for source in get_object_set(prev_world, Source) }
end

function get_type_to_score(action_to_score)
  return Dict{ # type => total score of actions on objs of this type
    type => sum(s for (a, s) in action_to_score if a.type == type)
    for type in unique(a.type for (a, _) in action_to_score)
  }
end

# use Gen `assess` function to get the PDF that `sample_properties`
# proposes `output` given `detected`
function prob_of_sampling(output, detected, scenelength)
  (logpdf, _) = assess(sample_properties,
    (output.type, detected, scenelength), choicemap(
      (output.type == noise ? :amp : :erb) => output.amp_or_erb,
      :onset => output.onset,
      :duration => output.duration
    ))
  return exp(logpdf)
end

involution function propose_detected_sound_involution()
  set backward[:do_birth] = |proposed[:do_birth]
  if proposed[:do_birth]
    set backward[:rev => :reverse_is_smart] = |proposed[:do_smart_birth]
    set backward[:action] = DeathAction(proposed[:objtype], proposed[:idx])
    set backward[:objtype] = proposed[:objtype]
    source = create Source with index proposed[:idx]

    set is_noise(source) = (proposed[:objtype] == noise)
    if proposed[:do_smart_birth]
      for prop in (:duration, :onset, is_noise ? :amplitude : :erb)
        set ($prop)(source) = proposed[:properties => prop]
      end
    else
      for prop in (:duration, :onset, is_noise ? :amplitude : :erb)
        regenerate ($prop)(source)
      end
    end
  end
  backward[:do_smart_birth] = proposed[:rev => :reverse_is_smart]
  source = delete Source with index proposed[:action].idx
  backward[:idx] = proposed[:action].idx
  if proposed[:rev => :reverse_is_smart]
    for prop in (:duration, :onset, is_noise ? :amplitude : :erb)
      backward[:properties => prop] = get(($prop)(source))
    end
    backward[:action] = proposed[:rev => :reverse_birth]
  else
    for prop in (:duration, :onset, is_noise ? :amplitude : :erb)
      backward[($prop)(source)] = regenerated
    end
  end
end
end
end
end

```

Figure 11: The **data-driven sound detection kernel** inference kernel used in the audio example. The kernel uses an image-filtering edge-detection algorithm *detect\_sounds* (not shown) to detect noises and tones in the observed sound waves which are not explained by the currently hypothesized audio sources. Heuristic scoring functions *birth\_score* and *death\_score* (not shown) are used to score the value of adding detected sounds to the latents, or deleting sounds in the current latents. Based on these scores, either a currently inferred sound is selected to be deleted (a *DeathAction*), or a detected sound is selected to be created, with some noise added to the detected sound parameters (a “birth” move). To ensure reversibility, death moves must select which detected sound will be sampled during the reversing birth move. If no compatible sounds will be detected during a reversing birth, the death move will specify that the reversing birth should be “dumb” (non-data-driven) and create a sound generated from the prior.



```

gen function split_at_detected_endpoints_proposal(prev_world)
  sources = get_object_set(prev_world, Source())
  tones = filter(s -> !get(prev_world, :is_noise, s), sources)
  noises = filter(s -> get(prev_world, :is_noise, s), sources)
  merge_possible = length(tones) > 1 || length(noises) > 1
  splitprior = merge_possible ? (length(sources) < 4 ? 0.5 : 0) : 1
  do_split ~ bernoulli(splitprior)
  if do_split
    to_split ~ uniform_choice(sources)
    idx1 ~ uniform_discrete(1, length(sources) + 1)
    idx2 ~ uniform_from_list({i for i=1:length(sources) + 1 if i != idx1})
    if get(prev_world, :is_noise, to_split)
      old_amp = get(prev_world, :amplitude, to_split)
      amp1 ~ normal(old_amp, 0.5); amp2 ~ normal(old_amp, 0.5)
    else
      old_freq = get(prev_world, :freq_in_erb, to_split)
      freq1 ~ normal(old_freq, 0.5); freq2 ~ normal(old_freq, 0.5)
    end
    {*} ~ smart_sample_durations(prev_world, to_split)
  else
    to_idx ~ uniform_discrete(1, length(sources) - 1)
    mergeable = (length(tones) > 1 && length(noises) > 1) ?
      union(noises, tones) : (length(tones) > 1 ? tones : noises)
    # sample 2 sources to merge, where the second's beginning
    # and ending is after the first's
    function is_before(x, y)
      startx = get(prev_world, :starttime, x)
      starty = get(prev_world, :starttime, y)
      durx = get(prev_world, :duration, x)
      dury = get(prev_world, :duration, y)
      return startx < starty && startx + durx < starty + dury
    end
    function is_not_last_mergeable(x)
      return length(filter(y -> is_before(y, x), mergeables)) > 1
    end
    not_last_mergeables = filter(is_not_last_mergeable, mergeables)
    merge1 ~ uniform_choice(not_last_mergeables)
    merge2 ~ uniform_choice(filter(x -> is_before(x, merge1), mergeables))
    if get(prev_world, :is_noise, merge1)
      amp1 = get(prev_world, :amplitude, merge1)
      amp2 = get(prev_world, :amplitude, merge2)
      amplitude ~ normal((amp1 + amp2)/2, 0.5)
    else
      freq1 = get(prev_world, :freq_in_erb, merge1)
      freq2 = get(prev_world, :freq_in_erb, merge2)
      freq ~ normal((freq1 + freq2)/2, 0.5)
    end
  end
end

# one of the 2 new sounds we split `to_split` to will start
# at the time when `to_split` begins; the other will end at the
# time `to_split` ends. This function samples how long each will be.
gen function smart_sample_durations(prev_world, to_split)
  # analyze the region of the spectrogram in the time and frequency
  # range of the sound we are splitting; use a Sobel edge-detection
  # filter to find breaks in the sound in this region
  likely_endpoints = get_detected_endpoints(prev_world, to_split)
  if likely_endpoints != nothing # if breakpoints were detected
    likely_start, likely_end = likely_endpoints/SCENE_LENGTH
    sound_start = get(prev_world, :starttime, to_split)
    sound_duration = get(prev_world, :duration, to_split)
    likely_duration1 = likely_start - sound_start
    likely_duration2 = sound_start + sound_duration - likely_end
    dur1 ~ normal(likely_duration1, 0.05)
    dur2 ~ normal(likely_duration2, 0.05)
  else
    dur1 ~ uniform(1, SCENE_LENGTH)
    dur2 ~ uniform(1, SCENE_LENGTH)
  end
end

involution function splitmerge_involution()
  backward[:do_split] = !proposed[:do_split]
  if proposed[:do_split]
    old = proposed[:to_split]
    (new1, new2) = split old to indices
      proposed[:idx1], proposed[:idx2]
    backward[:merge1] = new1
    backward[:merge2] = new2
    backward[:to_idx] = index(old)
    set is_noise(new1) = get(is_noise(old))
    set is_noise(new2) = get(is_noise(old))
    # set starttimes and durations of new sounds
    set starttime(new1) = get(starttime(old))
    set duration(new1) = proposed[:duration1]
    # the second new sound will have the same end as the old sound.
    # the model specifies in terms of starttime, duration, so we
    # must do some math to calculate the new startpoint:
    start2 = (get(starttime(old)) + get(duration(old))
      - proposed[:duration2])
    set starttime(new2) = start2
    set duration(new2) = proposed[:duration2]
  else
    if get(is_noise(old))
      set amplitude(new1) = proposed[:amp1]
      set amplitude(new2) = proposed[:amp2]
      backward[:amplitude] = get(amplitude(old))
    else
      set freq_in_erb(new1) = proposed[:freq1]
      set freq_in_erb(new2) = proposed[:freq2]
      backward[:freq] = get(freq_in_erb(old))
    end
  end
  old1, old2 = proposed[:merge1], proposed[:merge2]
  new = merge old1, old2 to index proposed[:to_idx]
  backward[:idx1] = index(old1)
  backward[:idx2] = index(old2)
  backward[:to_split] = new
  set is_noise(new) = get(is_noise(old1))
  set starttime(new) = get(starttime(old1))
  endpoint = get(starttime(old2)) + get(duration(old2))
  set duration(new) = endpoint - get(starttime(old1))
  backward[:duration1] = get(duration(old1))
  backward[:duration2] = get(duration(old2))
  if get(is_noise(old1))
    set amplitude(new) = proposed[:amplitude]
    backward[:amp1] = get(amplitude(old1))
    backward[:amp2] = get(amplitude(old2))
  else
    set freq_in_erb(new) = proposed[:freq]
    backward[:freq] = get(freq_in_erb(old1))
    backward[:freq] = get(freq_in_erb(old2))
  end
end
end
end

```

Figure 12: The **data-driven sound splitting kernel** either chooses to split one long sound into two short ones, or merge two short tones into a longer one. In the split move, we use a similar edge-detection algorithm to what was used in the **create detected sound** kernel to detect changepoints in the observed sound during the time at which the sound being split occurs. The kernel proposes an endpoint for the first new sound and a startpoint for the second new sound to be close to the detected changepoints.

proposing endpoints for the new sounds based on the analysis of the sounds. Figure 3(b) and Figure 8 show that the addition of this kernel substantially speeds up convergence. This is to be expected, since this kernel can integrate information from the data about another class of possible mistakes in inferred latents: when one sound has been mistakenly interpreted as two, or vice versa. We cycle both of the data-driven kernels with Gaussian drift kernels for the parameters of inferred audio sources.

#### A.2.1. THE AUTOMATED ACCEPTANCE PROBABILITY FOR THE **data-driven sound detection kernel**

To illustrate the calculations our system automates, we give the acceptance probability calculation performed when running the **data-driven sound detection kernel**, including calculating the Jacobian of the transformation.

**Model density.** We can factor the density of the model,  $p(x)$ , as  $p(x) = p(x_{\mathcal{L}})p(x_{\mathcal{O}}|x_{\mathcal{L}})$ , where  $x_{\mathcal{L}}$  is the collection of latent variable values in  $x$  (the number of underlying sound sources, and the properties of each one), and  $x_{\mathcal{O}}$  is the observed sound waves variable in  $x$ . For any  $x$  where each variable is within its support, the probability of the latents is:

$$\begin{aligned}
 p(x_{\mathcal{L}}) &= \frac{1}{4} \prod_{i=1}^{x_{N_S}} (p(x_{D,i}) \times p(x_{\tau,i}) \times p(x_{(\alpha|f),i}) \times p(x_{n,y_s})) \\
 &= \frac{1}{4} \prod_{i=1}^{x_{N_S}} \left( \left( \frac{1}{L-0.1} \right) \times \left( \frac{1}{L-0.1} \right) \times \left( \begin{cases} \mathcal{N}(x_{(\alpha|f),i}, 8, 10) & x_{n,i} = 1 \\ \frac{1}{24-0.4} & \text{otherwise} \end{cases} \right) \times 0.4^{x_{n,i}} (1-0.4)^{1-x_{n,i}} \right)
 \end{aligned}$$

Here,  $x_{N_S}$  is the number of sound sources in MSSSI  $x$  (the value of the number variable `number Source()`);  $x_{D,i}$  is the duration of the  $i$ th audio source in  $x$ ;  $x_{\tau,i}$  is the starttime of the  $i$ th audio source in  $x$ ;  $x_{n,i}$  is an indicator which is 1 if the  $i$ th source in  $x$  is a noise and 0 otherwise;  $x_{(\alpha|f),i}$  denotes the amplitude of the  $i$ th sound in  $x$  if it is a noise, and the frequency of this sound if it is a tone. We let  $L$  denote the length of the rendered audio tones (which is a global constant).

The probability of the observation given the latents is:

$$p(x_{\mathcal{O}} | x_{\mathcal{L}}) = \mathcal{N}(x_{\mathcal{O}}, r(x_{\mathcal{L}}), I_{|r(x_{\mathcal{L}})|})$$

where  $x_{\mathcal{O}}$  denotes the observed soundwaves in  $x$  (from the `observed_sound` property in Figure 7) represented as a vector,  $r(x_{\mathcal{L}})$  denotes the *rendering* of the latent audio sources in  $x$  into sound-waves, and  $I_k$  denotes the  $k$ -dimensional identity matrix.

**State-dependent proposal density.** The state-dependent proposal distribution  $q_x(y)$  is parameterized by a current MSSSI  $x$ . It is a distribution over finite dictionaries, implemented as a Gen program. Using standard techniques from probabilistic programming, Gen can automate the computation of densities of dictionaries sampled by the proposal:

$$q_x(y) = \begin{cases} \frac{(1 - \beta(x))\phi_n(x)_{y_n}\phi(x; y_n)_{y_i}(\beta_\rho(x, y_i)\theta_\rho(x, y_i; y_n)_{y_{s'}})^{y_\rho}(1 - \beta_\rho(x, y_i))^{1-y_\rho}}{\frac{\beta(x)\cdot(1-\beta_\sigma(x))}{x_{N_S}+1} \cdot \frac{(0.4\cdot\mathcal{N}(y_\alpha; 8, 10))^{y_n}\cdot(0.6\cdot(37-0.4)^{-1})^{1-y_n}}{(L-0.1)^2}} & y_b = 0 \\ & y_b = 1 \wedge y_\sigma = 0 \\ \frac{\beta(x)\cdot\beta_\sigma(x)}{x_{N_S}+1} \cdot \theta_n(x)_{y_n}\theta(x; y_n)_{y_s} \cdot \mathcal{N}_T \left( \begin{bmatrix} y_D \\ y_\tau \\ y_{v(y_n i)} \end{bmatrix}; \begin{bmatrix} R_D \\ R_\tau \\ R_{v(y_n)} \end{bmatrix}, \begin{bmatrix} D_D(x, y_s, y_n) \\ D_\tau(x, y_s, y_n) \\ D_{v(y_n)}(x, y_s, y_n) \end{bmatrix}, \begin{bmatrix} 0.1 & 0 & 0 \\ 0 & 0.1 & 0 \\ 0 & 0 & \sigma_{v(y_n)} \end{bmatrix} \right) & y_b = 1 \wedge y_\sigma = 1 \end{cases}$$

Here,  $x_{N_S}$  is the number of sound sources in MSSSI  $x$ , and  $y_b$  determines whether a new sound source will be proposed; if  $y_b = 1$ ,  $y_\sigma$  determines whether the new source will be proposed using the data-driven edge-detection proposal or from the prior. (In Figure 11,  $y_b$  corresponds to `do_birth`;  $y_\sigma$  to `do_smart_birth`.)  $y_n$  is an indicator for whether the source being created/deleted is white noise, as opposed to a tone ( $y_n$  is `objtype` in Figure 11 when  $y_b = 0 \vee y_\sigma = 1$ , and is proposed via `regenerate` when  $y_b = 1 \wedge y_\sigma = 0$ ). When  $y_\sigma = 1$ , an index  $y_s$  is computed determining *which* of the sounds from the image-filtering and edge-detection algorithm with the type indicated by  $y_n$  will be proposed (the index of sample `action` in Figure 11). Based on  $y_s$ , new sound onset and duration parameters  $y_\tau$  and  $y_D$ , as well as either a frequency  $y_f$  or an amplitude  $y_\alpha$ . When the data-driven proposal is not used ( $y_b = 1 \wedge y_\sigma = 0$ ),  $y_\alpha$  (or  $y_f$ ),  $y_n$ , and other parameters are proposed from the model's prior over sound sources. If  $y_b = 0$ , a death move is proposed, to delete sound  $y_i$ . The variable  $y_\rho$  is introduced to enable the implementation of an involution: it represents whether the move *reversing* a proposed death move will be data-driven. If  $y_\rho = 1$ ,  $y_{s'}$  represents which heuristically detected sound source with the type indicated by  $y_n$  will be used as the basis to repropose the deleted audio source, in the reverse proposal.  $\mathcal{N}_T(v; R, \mu, \Sigma)$  denotes the PDF of a multivariate truncated normal distribution with ranges  $R$ , means  $\mu$ , and covariance matrix  $\Sigma$ , evaluated at  $v$ .

The proposal is based on a number of functions, as well as the scene length parameter  $L$ , the standard deviations  $\sigma_\alpha = 4$  and  $\sigma_f = 1$ , and the ranges  $R_D = (0.1, 1)$ ,  $R_\tau = (0, L)$ ,  $R_f = (0.4, 24)$ , and  $R_\alpha = (-\infty, \infty)$ . These functions, and some additional functions needed for their definitions, are described below:

- $k(x, y_n)$  is the number of possible sounds to create of type  $y_n$  detected in state  $x$
- $S^B(x, i, y_n)$  is a score for the action of creating a sound close to the  $i$ th detected sound of type  $y_n$ , from a user-written scoring heuristic
- $S^D(x, i)$  is a score for the action of deleting the  $i$ th sound in the current world, from a user-written scoring heuristic (a different heuristic from  $S^B$ )
- $\pi_B(x) = \frac{\sum_{y_n=0}^1 \sum_{i=1}^{k(x, y_n)} S^B(x, i, y_n)}{\sum_{y_n=0}^1 \sum_{i=1}^{k(x, y_n)} S^B(x, i, y_n) + \sum_{i=1}^{x_{N_S}} S^D(x, i)}$  is the total score of all detected birth actions, divided by the total score of all detected birth and death actions

$$\beta(x) = \begin{cases} \frac{\pi_B(s) + \pi_B^* \pi_R}{1 + \pi_R} & x_{N_S} \in (0, 4) \\ 0 & x_{N_S} = 4 \\ 1 & x_{N_S} = 0 \end{cases}$$

is the probability of  $y_b = 1$ , given  $x$ . Hyperparameter  $\pi_B^* = 0.5$  is a prior probability of doing a birth move, not factoring in the scores assigned to detected actions, and hyperparameter  $\pi_R = 0.1$  is the probability of making a decision to do a birth vs death move without considering the scores.

- $\beta_\sigma(x) = 1 - 0.05^{[\theta(x)] > 0]}$  is the probability of  $y_\sigma = 1$  given that  $y_b = 1$  for MSSI  $x$ .
- $\beta_\rho(x, i)$  is 0, if removing sound  $i$  from  $x$  and running image filtering and edge detection does not find any sounds of the same type as sound  $i$  in  $x$ , and 0.95 otherwise.
- $D_D(x, i, y_n)$  is the predicted duration of the  $i^{\text{th}}$  detected sound source of type  $y_n$  from the image filtering and edge detection algorithm applied to  $x$ . Similarly,  $D_\tau$  gives the predicted onset,  $D_\alpha$  the predicted amplitude (if applicable), and  $D_f$  the predicted frequency (if applicable).
- $v(y_n)$  is the symbol  $\alpha$  (amplitude) if  $y_n = 1$  and the symbol  $f$  (frequency) if  $y_n = 0$ .
- $\theta_n(x) = \{0 \mapsto \frac{\sum_{i=1}^{k(x,0)} S^B(x,i,0)}{\sum_{i=1}^{k(x,0)} S^B(x,i,0) + \sum_{i=1}^{k(x,1)} S^B(x,i,1)}, 1 \mapsto \frac{\sum_{i=1}^{k(x,1)} S^B(x,i,1)}{\sum_{i=1}^{k(x,0)} S^B(x,i,0) + \sum_{i=1}^{k(x,1)} S^B(x,i,1)}\}$  is a dictionary from values of  $y_n$  to the probability of choosing to birth an object of type  $y_n$  based on the heuristically-assigned  $S^B$  scores.
- $\phi_n(x) = \{0 \mapsto \frac{\sum_{i \in \mathbb{N} \cap (0, x_{N_S}): x_{n,i}=0} S^D(x,i)}{\sum_{i=1}^{x_{N_S}} S^D(x,i)}, 1 \mapsto \frac{\sum_{i \in \mathbb{N} \cap (0, x_{N_S}): x_{n,i}=1} S^D(x,i)}{\sum_{i=1}^{x_{N_S}} S^D(x,i)}\}$  gives the probability of deleting a noise vs deleting a tone.
- $\theta(x; y_n)$  is a vector-valued function, taking values in the  $(k(x, y_n) - 1)$ -simplex. It represents a distribution over more or less promising detections based on the  $S^B$  scores. In particular,  $\theta(x; y_n)_j = \frac{S^B(x, j, y_n)}{\sum_{i=1}^{k(x, y_n)} S^B(x, i, y_n)}$ .
- $\phi(x; y_n)$  is a vector-valued function representing a distribution over which existing audio source of type  $y_n$  to delete. In particular,  $\phi(x; y_n)_j = \frac{S^D(x, j)}{\sum_{i \in \mathbb{N} \cap (0, x_{N_S}): x_{n,i}=0} S^D(x, i)}$ .
- $\theta_\rho(x, i; y_n)$  is a vector-valued function, taking values in the  $(k_\rho(x, i, y_n) - 1)$  simplex, where  $k_\rho(x, i, y_n)$  is the number of detected sounds of type  $y_n$  from the edge detection routine, applied to  $x$  without sound source  $i$ . In particular,

$$\theta_\rho(x, i; y_n)_j \propto \mathcal{N}_T \left( \begin{bmatrix} x_{D,i} \\ x_{\tau,i} \\ x_{v(y_n),i} \end{bmatrix}; \begin{bmatrix} R_D \\ R_\tau \\ R_{v(y_n)} \end{bmatrix}, \begin{bmatrix} D_D^\rho(x, i, j, y_n) \\ D_\tau^\rho(x, i, j, y_n) \\ D_{v(y_n)}^\rho(x, i, j, y_n) \end{bmatrix}, \begin{bmatrix} 0.1 & 0 & 0 \\ 0 & 0.1 & 0 \\ 0 & 0 & \sigma_{v(y_n)} \end{bmatrix} \right)$$

where  $D_D^\rho$ ,  $D_\tau^\rho$ ,  $D_\alpha^\rho$ , and  $D_f^\rho$  are functions on  $(x, i, j, y_n)$  giving the duration, onset, and amplitude/frequency of the  $j$ th sound of type  $y_n$  detected via the detector applied to  $x$  without source  $i$ .

**The Jacobian.** The **create detected sound** kernel always either creates or deletes a sound; we show the Jacobian of the matrix for a move which creates a sound (the transpose of which is the Jacobian of the matrix for the corresponding deletion). The transition  $(x, y) \rightarrow (x', y')$  during a creation copies over every continuous value in  $x$  and  $y$  into  $x'$ . These values come from the properties **starttime** ( $\tau$ ) and **duration** ( $D$ ) of each source, and the **amplitude** ( $D$ ) property of each noise or **frequency** ( $f$ ) property of each tone. Depending on the index of the newly-created sound, the indices of each object may change, as is described in . In particular, when we birth a new source with index  $i$ , we must increment the index of each existing source with index  $\geq i$ . We define the following identifiers for whether the object at index  $a$  is moved to index  $b$ , and whether the birth index  $i$  equals a given index  $a$ :

$$C_{a \rightarrow b} = \begin{cases} 1 & (i > a \wedge a = b) \vee (i \leq a \wedge b = a + 1) \\ 0 & \text{otherwise} \end{cases}$$

$$B_a = \begin{cases} 1 & i = a \\ 0 & \text{otherwise} \end{cases}$$

Using these definitions, we can write the Jacobian of the transformation which copies starttime ( $\tau$ ), duration ( $D$ ), and amplitude/frequency ( $\alpha|f$ ) values for each object as:

$$J = \begin{matrix} & x_{(\alpha|f),1} & x_{\tau,1} & x_{D,1} & x_{(\alpha|f),2} & x_{\tau,2} & x_{D,2} & \dots & x_{D,N_S+1} \\ \begin{matrix} y_{(\alpha|f)} \\ y_{\tau} \\ y_D \\ x_{(\alpha|f),1} \\ x_{\tau,1} \\ x_{D,1} \\ x_{(\alpha|f),2} \\ x_{\tau,2} \\ x_{D,2} \\ \vdots \\ x_{D,N_S} \end{matrix} & \left( \begin{matrix} B_1 & 0 & 0 & B_2 & 0 & 0 & \dots & 0 \\ 0 & B_1 & 0 & 0 & B_2 & 0 & \dots & 0 \\ 0 & 0 & B_1 & 0 & 0 & B_2 & \dots & B_{N_S+1} \\ C_{1 \rightarrow 1} & 0 & 0 & C_{1 \rightarrow 2} & 0 & 0 & \dots & 0 \\ 0 & C_{1 \rightarrow 1} & 0 & 0 & C_{1 \rightarrow 2} & 0 & \dots & 0 \\ 0 & 0 & C_{1 \rightarrow 1} & 0 & 0 & C_{1 \rightarrow 2} & \dots & C_{1 \rightarrow N_S+1} \\ C_{2 \rightarrow 1} & 0 & 0 & C_{2 \rightarrow 2} & 0 & 0 & \dots & 0 \\ 0 & C_{2 \rightarrow 1} & 0 & 0 & C_{2 \rightarrow 2} & 0 & \dots & 0 \\ 0 & 0 & C_{2 \rightarrow 1} & 0 & 0 & C_{2 \rightarrow 2} & \dots & C_{2 \rightarrow N_S+1} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \\ 0 & 0 & C_{2 \rightarrow N_S} & 0 & 0 & C_{2 \rightarrow N_S} & \dots & C_{N_S \rightarrow N_S+1} \end{matrix} \right) \end{matrix}$$

where

$$J_{a,b} = \begin{cases} 0 & a \pmod{3} \neq b \pmod{3} \\ B_b & a \leq 3 \\ C_{a \rightarrow b} & a > 3 \end{cases}$$

Since this is a permutation matrix, it always has determinant  $-1$  or  $1$ . Our system is capable of computing this Jacobian using automatic differentiation and discovering this. As an optimization, in practice, our system can identify that this transformation only copies values, and thus will have a permutation-matrix Jacobian which it does not need to take the determinant of, since it will certainly have magnitude 1, as we describe in the **Jacobian sparsity** paragraph in Section D. For an description of an inference kernel we have implemented which has a Jacobian whose determinant's magnitude is (usually) not 1, see Section A.1.

**Acceptance probability.** When using the involutive MCMC kernel, our system automatically calculates  $q_x(y)$ ,  $q_{x'}(y')$ ,  $p(x)$ ,  $p(x')$ , and  $J$ , to compute the acceptance probability of the move  $(x, y) \rightarrow (x', y')$ ,

$$\frac{p(x')}{p(x)} \times \frac{q_{x'}(y')}{q_x(y)} \times |\det J|$$

These terms are all equivalent to the above expressions for  $q_x(y)$ ,  $p(x)$ , and  $J$  (sometimes substituting  $x'$  for  $x$  or  $y'$  for  $y$ ). In practice, the  $q$  terms are automatically calculated using the given expressions, while computational optimizations let us avoid calculating the determinant of the Jacobian matrix (which is just a permutation matrix), and let us avoid fully calculating  $p(x')$  and  $p(x)$ , instead directly calculating

$$\begin{aligned} \frac{p(x')}{p(x)} &= \frac{p(x'_{\mathcal{O}}|x'_{\mathcal{L}}) \times \frac{1}{4} \prod_{i=1}^{x'_{NS}} p(x'_{D,i}) \times p(x'_{\tau,i}) \times p(x'_{(\alpha|f),i}) \times p(x'_{n,i})}{p(x_{\mathcal{O}}|x_{\mathcal{L}}) \times \frac{1}{4} \prod_{i=1}^{x_{NS}} p(x_{D,i}) \times p(x_{\tau,i}) \times p(x_{(\alpha|f),i}) \times p(x_{n,i})} \\ &= \frac{p(x'_{\mathcal{O}}|x'_{\mathcal{L}})}{p(x_{\mathcal{O}}|x_{\mathcal{L}})} \times \begin{cases} p(x'_{D,y_s}) \times p(x'_{\tau,y_s}) \times p(x'_{(\alpha|f),y_s}) \times p(x'_{n,y_s}) & y_b = 1 \\ \frac{1}{p(x_{D,y_s}) \times p(x_{\tau,y_s}) \times p(x_{(\alpha|f),y_s}) \times p(x_{n,y_s})} & y_b = 0 \end{cases} \end{aligned}$$

where, as above,  $y_b$  is 1 when  $y$  proposes to create a new source and 0 otherwise, and  $y_s$  is the index of the source being created/deleted. We can avoid evaluating the product expressions because the majority of terms in the numerator and denominator cancel out, since all audio sources except the one being created/deleted are unaffected by the update.



## Appendix B. Open-Universe Models and Contingent Bayesian Networks

In this section, we give a more complete definition of open-universe probabilistic models and contingent Bayesian networks (outlined in Section 2).

Open-universe probabilistic models are models of relational domains in which the number of objects, and their link structure, is not known a priori. The usual representations for structured probability distributions—e.g., probabilistic graphical models—are insufficient for specifying OUPMs, which typically define an unbounded number of random variables, with conditional independence relationships that are themselves random. OUPMs can be specified formally, however, as programs in sufficiently expressive probabilistic programming languages (Milch et al., 2005a), or as *contingent Bayesian networks* (CBNs) with infinitely many nodes (Milch et al., 2005b). Here, we review the standard definition of OUPMs as CBNs, and extend it to allow for objects with real-valued random properties.

**Open-Universe Probabilistic Models.** An *open universe probabilistic model* is a tuple  $(\mathcal{T}, \mathcal{N}, \mathcal{P}, \mathcal{C})$ , where

- $\mathcal{T}$  is a finite<sup>3</sup> set of *object types* over which the model is defined (e.g., *Person*).
- $\mathcal{N}$  is a finite<sup>3</sup> set of *object origin declarations*  $(\tau, (\tau_1, \dots, \tau_n))$ , with  $\tau, \tau_j \in \mathcal{T}$ . An object origin declaration specifies that a tuple of objects, of types  $(\tau_1, \dots, \tau_n)$  respectively, can be the *origin* of a separate object of type  $\tau$ . For example, in a model of genealogy, we might declare  $(\textit{Person}, (\textit{Person}, \textit{Person})) \in \mathcal{N}$ , encoding that a pair of people can be the origin (i.e., parents) of other people. A model’s object origin declarations determines sets of *possible objects*  $\mathcal{O}_\tau(\mathcal{N})$  for each type  $\tau$ , which are finite expressions of the form  $(\tau, (o_1, \dots, o_n), i)$ , where  $(\tau, (\tau_1, \dots, \tau_n)) \in \mathcal{N}$ ,  $i \in \mathbb{N}$ , and  $o_j \in \mathcal{O}_{\tau_j}(\mathcal{N})$  for each  $j$ . We call  $(\tau, (o_1, \dots, o_n), i)$  the  $i^{\text{th}}$  object of type  $\tau$  with origin  $(o_1, \dots, o_n)$ .
- $\mathcal{P}$  is a countable set of typed *object property declarations*  $(P, (\tau_1, \dots, \tau_n))$ , with  $\tau_j \in \mathcal{T}$ , together with value domains  $\text{dom}(P)$  for each property. A property represents data concerning zero or more objects, of types  $(\tau_1, \dots, \tau_n)$  respectively. For example, a model might declare a property representing the number of flights between two places.
- $\mathcal{C}$  is a *continuously-extended contingent Bayesian network* on the set of *possible variables* for  $(\mathcal{T}, \mathcal{N}, \mathcal{P})$ ,  $\{(x, (o_1, \dots, o_n)) \mid (x, (\tau_1, \dots, \tau_n)) \in \mathcal{P} \cup \mathcal{N} \wedge \forall j, o_j \in \mathcal{O}_{\tau_j}(\mathcal{N})\}$ . For every origin declaration  $(\tau, (\tau_1, \dots, \tau_n)) \in \mathcal{N}$ , and every tuple  $(o_1, \dots, o_n) \in \bigotimes_j \mathcal{O}_{\tau_j}(\mathcal{N})$  of possible objects of the correct types, there is a possible *number variable*  $(\tau, (o_1, \dots, o_n))$ , with domain  $\mathbb{Z}_{>0}$ , representing the number of objects of type  $\tau$  that have objects  $(o_1, \dots, o_n)$  as their origin. For every property declaration  $(P, (\tau_1, \dots, \tau_n))$  and object tuple  $(o_1, \dots, o_n)$  of the correct types, there is a possible *property variable* specifying the value of the property  $P$  (within  $\text{dom}(P)$ ) of the object tuple  $(o_1, \dots, o_n)$ . The contingent Bayesian network itself defines a distribution over assignments to these infinitely many possible variables, as we describe in the following section.

3. The set of types and the set of object origin declarations can be countable without compromising any results, but existing formal languages for defining OUPMs, including ours, do not provide a way to declare models where either of these are infinite. It would be interesting to consider whether any natural applications would require this flexibility, and how languages might be designed to expose it.

**Contingent Bayesian Networks.** For a set of possible variables  $\mathcal{V}$ , let  $\Omega = \bigotimes_{v \in \mathcal{V}} \text{dom}(v)$  be the space of all possible assignments of the (possibly infinitely many) variables to values within their domains. When the domains are discrete, we can use a *contingent Bayesian network* (CBN) (Milch et al., 2005b) to define a distribution over such assignments. Analogously to a Bayesian network on a finite number of variables, a CBN on an infinite set of variables  $\mathcal{V}$  specifies an infinite directed graph with  $\mathcal{V}$  as the nodes. But whereas traditional Bayesian networks fix the dependency structure, in a CBN the edges are *contingent*: each edge  $u \rightarrow v$  is labeled with an event  $\Omega_{u \rightarrow v} \subseteq \Omega$  within which the edge is *active*. As in Bayesian networks, to each variable  $v$  we associate a conditional probability distribution (CPD)  $p_v$ . But this conditional probability distribution is allowed to depend only on the *active* parents of  $v$ , which can be formalized as follows: for each variable  $v \in \mathcal{V}$ , the contingent Bayesian network associates a partition  $\Lambda_v$  of the space  $\Omega$ , such that if  $\omega_1$  and  $\omega_2$  differ only in their assignment to  $u$ , and the edge  $u \rightarrow v$  either does not exist or is *inactive* in  $\omega_1$ , then  $\omega_1$  and  $\omega_2$  belong to the same component of  $\Lambda_v$ . Then  $p_v(\cdot \mid \omega \in \lambda)$  is allowed to depend on *which component*  $\lambda$  of the partition  $\omega$  is in, but nothing else, preventing dependencies on the values of inactive parents. We write  $Pa_\omega(v)$  for the active parents of the variable  $v$  in state  $\omega$ . These conditional probabilities together define a probability measure  $\Pi$  on the space  $\Omega$ , equipped with a non-standard  $\sigma$ -algebra; we refer the reader to Milch et al. (2005b) for a detailed treatment.

**Minimal Self-Supporting Instantiations.** In open-universe models, possible worlds  $\omega \in \Omega$  are generally infinite, because they specify assignments to all of the possible variables  $v \in \mathcal{V}$ . But for inference, we must work with finite representations. If we can partition  $\Omega$  into disjoint equivalence classes  $\Omega_i$ , then we can work with the equivalence classes, using the distribution  $p(\Omega_i) = \Pi(\Omega_i)$ . One useful strategy for partitioning  $\Omega$  is to consider the *minimal self-supporting instantiations* of a finite set of variables  $U \subseteq \mathcal{V}$ . Each component  $w$  of the partition is labeled by a set of variables  $\text{vars}(w) \subseteq \mathcal{V}$ , with  $U \subseteq \text{vars}(w)$ , as well as an assignment  $w[v] \in \text{dom}(v)$  for each  $v \in \text{vars}(w)$ . The component  $w$  contains all full possible worlds  $\omega \in \Omega$  that agree with  $w$  on the values of each  $v \in \text{vars}(w)$ . To ensure that the components  $w$  are indeed disjoint, we only include  $w$  as a component if: (1) it is *self-supporting*, meaning that for all  $\omega \in w$ , if  $v \in \text{vars}(w)$  and  $u \in Pa_\omega(v)$ , then  $u \in \text{vars}(w)$ ; and (2) it is *minimal*, meaning that removing any variable  $v \notin U$  from  $\text{vars}(w)$  would cause it to fail to be self-supporting. Together, these criteria ensure that for any  $U \subseteq \mathcal{V}$ , the components  $w$  are disjoint, and so  $p(w) = \Pi(w)$  is a valid probability distribution over minimal self-supporting instantiations  $w$ . Conveniently, the probability of each component  $w$  is a simple product of the conditional distributions  $p_v$ , as in a typical Bayesian network:

$$p(w) = \Pi(w) = \prod p_v(w[v] \mid \{w[u] \mid u \in Pa_w(v)\}). \quad (4)$$

### B.1. Contingent Bayesian Networks with Continuous Variables

We now extend the CBN framework to support models with some continuous object property domains, so long as no discrete variables count continuous variables as parents.<sup>4</sup> (Wu

4. See the following section for a similar but more involved extension to a more general case, where discrete variables can depend on continuous variables.

et al. (2018) give a semantics for open-universe models that interprets them as infinite product measures over assignments to all basic random variables. This allows them to interpret OUPMs with continuous and discrete random variables, but does not serve our purposes: it does not yield a formula for the *density* of an MSSI in an open-universe model, nor can it be used to prove that incremental computation strategies for fast MCMC are correct, as—unlike in CBNs—it does not account for contingent dependency structure among variables.)

As before, let  $\mathcal{V}$  be a collection of variables, but now suppose some variables  $\mathcal{V}_c \subseteq \mathcal{V}$  have continuous domains. Then define the surrogate discrete domain  $\text{dom}_D(v)$  of a variable  $v$  to be  $\text{dom}(v)$  if  $v \notin \mathcal{V}_c$ , and  $\{\star\}$  otherwise. A *continuously extended* CBN defines (1) an ordinary CBN on the variables  $\mathcal{V}$ , treating them as having surrogate discrete domains  $\text{dom}_D$ , and (2) conditional probability densities  $p_v(\cdot \mid \{w[u] \mid u \in Pa_w(v)\})$  for each continuous variable  $v \in \mathcal{V}_c$ . The CBN induces a distribution  $\Pi$  on minimal self-supporting worlds, which we modify to define a probability measure on the space  $\mathcal{W}$  of *continuously extended minimal self-supporting worlds*: pairs  $(w_d, w_c)$  of a minimal self-supporting world  $w_d$  for the discrete CBN, plus a real-valued assignment  $w_c$  to all variables in  $\text{vars}(w_d) \cap \mathcal{V}_c$ . We let  $w[v]$  denote  $w_c[v]$  if  $v \in \mathcal{V}_c$  and  $w_d[v]$  otherwise. The distribution has density

$$p(w) = \Pi(w_d) \times \prod_{v \in \mathcal{V}_c \cap \text{vars}(w)} p_v(w_c[v] \mid \{w[u] \mid u \in Pa_w(v)\}) \quad (5)$$

with respect to the reference measure

$$\mu(E) = \sum_w \lambda^{|\mathcal{V}_c|}(\{w'_c \mid w' \in E \wedge \text{vars}(w') = \text{vars}(w) \wedge \forall v \in \text{vars}(w') \setminus \mathcal{V}_c, w'[v] = w[v]\}),$$

where  $\lambda^k$  denotes the  $k$ -dimensional Lebesgue measure. Using Equation 4 for the first term of  $p(w)$ , we see that in the continuous case, we still have the ordinary Bayes net factorization  $p(w) = \prod_{v \in \text{vars}(w)} p_v(w[v] \mid \{w[u] \mid u \in Pa_w(v)\})$ .

### B.1.1. CONTINUOUSLY-EXTENDED CBNs WITH CONTINUOUS PARENTS OF DISCRETE VARIABLES

In the previous section, we define continuously-extended contingent Bayesian networks with the restriction that no discrete variables have continuous variables as parents. We now give a more flexible development.

Let  $\mathcal{V}$  be a collection of variables, where  $\mathcal{V}_c \subseteq \mathcal{V}$  have continuous domains. Define surrogate discrete domain  $\text{dom}_D(v)$  to be  $\text{dom}(v)$  if  $v \notin \mathcal{V}_c$ , and  $\{\star\}$  otherwise. Now define a continuously extended CBN to comprise: (1) an ordinary CBN on the variables  $\mathcal{V}$ , treating them as having surrogate discrete domains  $\text{dom}_D$ , and (2) a conditional *joint* probability density  $\mathcal{P}_c(w_c \mid w)$  on the values of the continuous variables in a minimal self-supporting instantiation  $w$  given the values of all  $w$ 's discrete variables. This CBN induces distribution  $\Pi$  on minimal self-supporting worlds; as before, we modify this to define the space  $\mathcal{W}$  of *continuously extended minimal self-supporting worlds*, consisting of a minimal self-supporting world  $w$  of the CBN, and an assignment  $w[v]$  to each continuous variable in  $\text{vars}(v) \cap \mathcal{V}_c$ . We use  $w_c$  to denote the assignment of values  $w[v]$  to the continuous variables  $v \in \mathcal{V}_c \cap \text{vars}(w)$  in  $w$ . The distribution then has density

$$p(w) = \Pi(\{\omega \mid \forall v \in \text{vars}(w) \setminus \mathcal{V}_c, \omega[v] = w[v]\}) \times \mathcal{P}_c(w_c \mid w)$$

with respect to the same reference measure from Section 2:

$$\mu(E) = \sum_w \lambda^{|w_c|}(\{w'_c \mid w' \in E \wedge \text{vars}(w') = \text{vars}(w) \wedge \forall v \in \text{vars}(w') \setminus \mathcal{V}_c, w'[v] = w[v]\}),$$

where  $\lambda^k$  denotes the  $k$ -dimensional Lebesgue measure. We note that the density we give ultimately factorizes, as before, as a marginal probability for the discrete variables and a conditional density for the continuous variables. (In principle, this alone is not a restriction, because any distribution over minimal self-supporting instantiations can be factored in this way.) But in the previous formulation (from Section 2), this factorization was embedded into the discrete CBN itself, affecting which worlds  $w$  counted as minimal and self-supporting in the first place. The key difference between this section's formulation and the one in Section 2 is that here we divorce the factorization of the continuous density from the CBN's structure; the edges involving continuous variables in the CBN serve only to determine which continuous variables are required in a minimal, self-supporting instantiation for an observed set of variables  $U$ .

## Appendix C. Modeling and Inference Programming Languages for OUPMs

In this section, we briefly outline our modeling language for OUPMs, and our DSL for writing involutions for involutive MCMC on OUPMs.

**OUPM Modeling Language** We have developed a new modeling language that is inspired by BLOG (Milch, 2006) for use in the Gen probabilistic programming system; for examples, refer to Figure 1, Figure 5, or Figure 7. Our language extends one of Gen’s built-in modeling languages in the following ways. We add `type` statements for declaring object types, `property` declarations for adding properties, and `number` declarations for adding possible origins. When writing world models with our system, users must define a special property  $O()$  with respect to which the model will take minimal self-supporting instantiations; they do this using an `observation_model` statement. Property and number declarations are attached to *function bodies* which together define the contingent Bayesian network.

In particular, a property statement declares the existence of a property  $P(\tau_1, \dots, \tau_n)$ , and defines a probabilistic generative process for the property’s value. A number statement  $(\tau, (\tau_1, \dots, \tau_n))$  declares that  $(\tau_1, \dots, \tau_n)$  is a permissible origin signature for objects of type  $\tau$ , and declares the distribution over how many objects of type  $\tau$  there are with each origin matching  $(\tau_1, \dots, \tau_n)$ . (So each property/number statement simultaneously serves as an object property declaration/object origin declaration, and declares the distribution over the corresponding property/number variables in the OUPM’s CBN.) The bodies of number and property declarations are specified in an imperative probabilistic programming language that induces distributions over finite dictionaries, called traces. The language uses Julia syntax for defining functions, but may sample random variable values using  $x \sim d$  statements. A program in the language defines a distribution over traces mapping each key  $x$  to the value sampled from the corresponding  $d$  in an execution of the function. The language is a built-in modeling language in Gen that was augmented with support for `get( $P(o_1, \dots, o_n)$ )` statements to access the value of  $P(o_1, \dots, o_n)$  in the world; this introduces an edge in the contingent Bayes net representing the model.

Represented as imperative probabilistic programs, number variables must always have 0 or 1 keys in their traces, corresponding to either deterministically returning a number, or sampling a single random number. Properties may have more than 1 key in their trace (so long as the set of possible keys is countable), but when they do, the the property statement  $P$  must be understood as defining a separate object property declaration  $P \Rightarrow k$  for each  $k$  which can appear as a key in the trace. For simplicity, in what follows, we will assume that each  $P(o_1, \dots, o_n)$  is itself a variable with only 0 or 1 keys in its probabilistic program’s traces; however, the algorithms we give also hold when there can be multiple keys. We represent MSSIs as sets of pairs  $(P(o_1, \dots, o_n), t)$ , where  $P$  is a number or property probabilistic program whose type signature matches  $(o_1, \dots, o_n)$ , and  $t$  is a trace for this variable.

**OUPM Involution Programming Language** We implement a new OUPM involution programming language to implement involutions on the space of pairs  $(x, y)$ , where  $x$  is a world represented as described above, and  $y$  is a trace from a imperative probabilistic

program. Involutions are written as imperative functions which may call special commands to modify state  $x$  in order to produce  $x'$ , to access values in  $x$  and  $y$ , and to specify the new values for  $y'$ . Specifically, the DSL provides support for the following.

- **Querying the current model state ( $x$ ).** Users may query the current model state using the syntax `get( $P(o_1, \dots, o_n)[a]$ )`, which retrieves the sample at address  $a$  of the trace for the property  $P(o_1, \dots, o_n)$ . Similarly, `get_number( $\tau(o_1, \dots, o_n)$ )` retrieves the number variable value in  $x$  defining how many objects of type  $\tau$  exist with origin  $(o_1, \dots, o_n)$ . Several higher-level commands are also provided to access higher-level discrete information about the current state, including `get_object_set( $\tau$ )`, which gets the set of all existing objects of type  $\tau$  in  $x$ .
- **Querying the auxiliary randomness ( $y$ ).** The syntax `proposed[ $a$ ]` returns  $y[a]$ , the value at address  $a$  in the random trace  $y$  of the proposal program  $\mathcal{Q}$ .
- **Proposing the new model state ( $x'$ ).** The new model state is proposed via commands that imperatively modify the old state  $x$ . They can do this either by changing the *structure* of  $x$  (by adding or removing objects, and changing the origins or indices of objects), or by writing new values to object properties.

The commands **create**, **delete**, **split**, **merge**, and **change** add, remove, and change the relationships of objects in the possible world, updating the values of number variables as necessary. The behavior of each is detailed below, and illustrated and described in Figure 13. Commands that remove objects have the effect of deleting any traces for variables  $G(o_1, \dots, o_i, \dots, o_n)$  of an object tuple including an  $o_i$  slated for deletion. If the deleted object is the origin of other objects (e.g., if an **Event** with **Detections** is deleted), those objects are not automatically deleted: users may either manually remove them or change their origins using **change**.

The command `set  $G(o_1, \dots, o_n)[a] = v$`  sets the sample at address  $a$  of a trace for  $G(o_1, \dots, o_n)$  to  $v$  in the proposed state  $x'$ .

- **Transforming the auxiliary randomness ( $y'$ ).** The user’s code is responsible for computing a new trace  $y'$  of  $\mathcal{Q}$  that makes  $f$  into an involution. Intuitively,  $y'$  must specify the proposal choices that would *reverse* the current move; in Figure 2, the forward move that posits a new **Event** to explain a false positive is reversed by a move that deletes the newly created **Event** and reclassifies its sole detection as a false positive. Users can write to  $y'$  using the syntax `backward[ $a$ ] =  $v$` .

### C.1. Object manipulation moves

The involution programming language includes commands to create, delete, split, merge, and change the origins of objects in the world. It is clear that these moves should increment or decrements number variables in the world to change the number of objects hypothesized to exist. However, to enable us to delete any object with any index, or create an object at any index (needed for reversibility), we need some scheme for reindexing the other existing objects. (For example, when creating a new object 1, the system must “make room” for the new object 1, for example by incrementing the index of every existing object.) We describe the reindexing scheme we use.



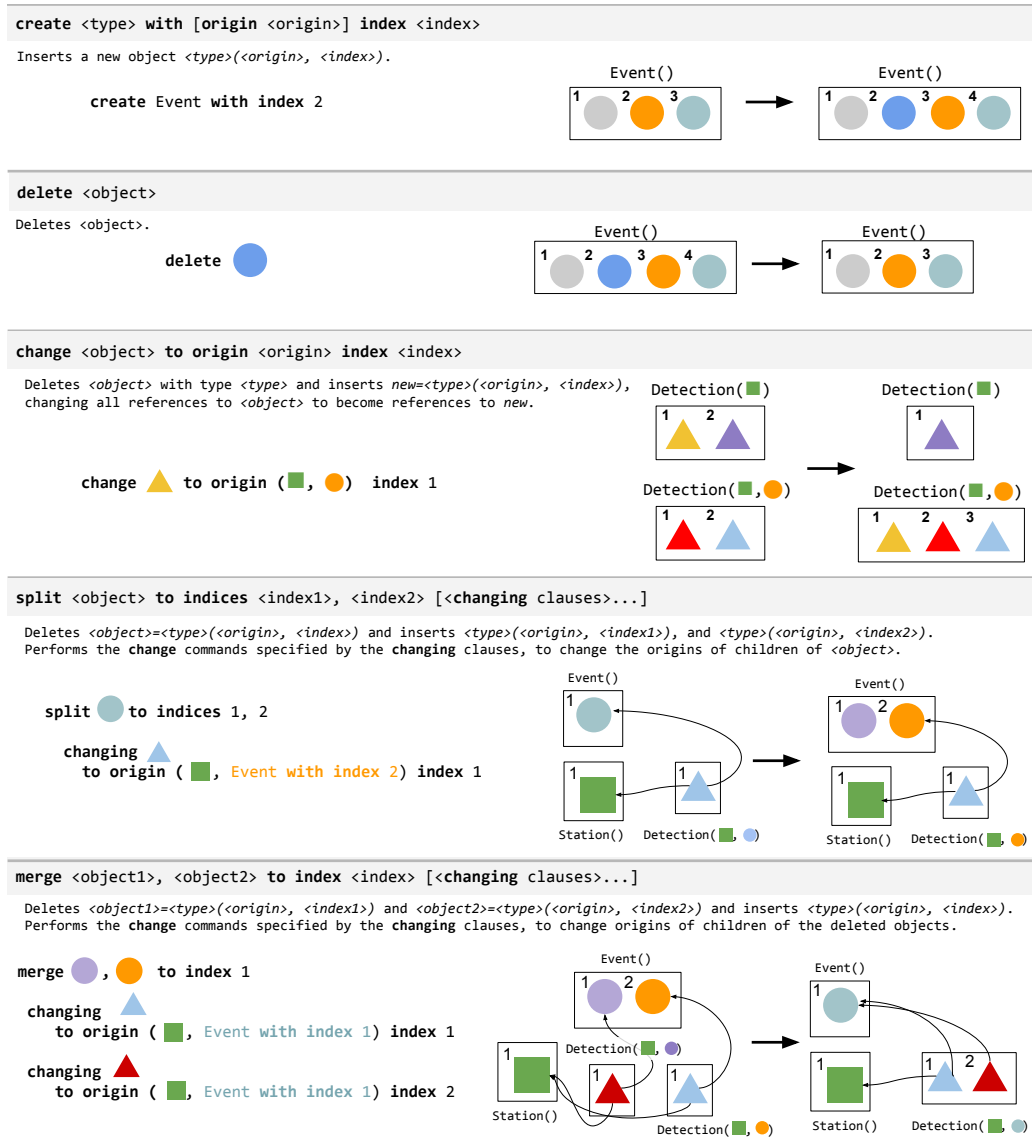


Figure 13: The object manipulation moves, with examples illustrating how they change the space of objects in the seismic monitoring example. To construct references objects, for use in the commands, illustrated with the shapes in this diagram, users use the syntax <type> with origin <origin> index <index>.

**Reindexing.** For any given origin tuple of objects, **origin**, a world contains  $N_T^{\text{origin}}$  objects of type  $T$ :  $\{T(\text{origin}, 1), \dots, T(\text{origin}, N_T^{\text{origin}})\}$ . Deleting the object  $T(\text{origin}, i)$  corresponds to a “vector deletion” which decrements the value of  $N_T^{\text{origin}}$  and decrements the index of every object  $T(\text{origin}, j)$  for  $j > i$ . Likewise, adding an object  $T(\text{origin}, i)$  is a “vector insertion” which increments  $N_T^{\text{origin}}$  and increments the index of every object  $T(\text{origin}, j)$  for  $j \geq i$ . In a world model, “changing the index” of an object  $T(\text{origin}, j)$  corresponds to updating the world so that the trace for any property or number variable  $P[o_1, \dots, T(\text{origin}, j), \dots, o_n]$  becomes the trace for  $P[o_1, \dots, T(\text{origin}, j + 1), \dots, o_n]$  or

$P[o_1, \dots, T(\text{origin}, j - 1), \dots, o_n]$ . Likewise, the update must change any object which has  $T(\text{origin}, j)$  in its origin to instead have  $T(\text{origin}, j + 1)$  or  $T(\text{origin}, j - 1)$  in its origin. While in a naive implementation, simple reindexing operations require copying many traces and greatly perturbing the model state, which causes extreme inefficiency when implementing split/merge and birth/death moves in PPLs like vanilla Gen or Stochaskell, the use of explicit open universe modeling constructs encodes enough structure for us to implement these moves very efficiently without copying data, using the object aliasing technique presented in Appendix D.

## Appendix D. Automatically Implementing Inference Kernels Efficiently

In this section, we describe in more detail the technique by which we automate the implementation of inference kernels, and the optimizations we use to produce implementations performant enough to scale to real-world data sizes. This technique extends upon an approach for automating involutive MCMC over traces of probabilistic programs described in Cusumano-Towner et al. (2020). In particular, we are able to exploit sparsity in the calculated Jacobian to avoid calculating the determinant of a large matrix, and exploit the conditional-independence relationships encoded in the contingent Bayesian network underlying the OUPM to efficiently calculate the state ratio  $p(x')/p(x)$ . The performance enabled by the latter is one of the key benefits of introducing first-class support for involutive MCMC on OUPMs.

To implement the inference kernel specified by a given proposal  $q$  and involution  $f$ , we must be able to sample and score auxiliary randomness  $y \sim q_x$ , produce transitions  $(x, y) \rightarrow (x', y')$ , and calculate acceptance probabilities by calculating  $q_x(y)$ ,  $q_{x'}(y')$ ,  $\frac{p(x')}{p(x)}$ , and  $|\det \frac{\partial(x', y')}{\partial(x, y)}|$ <sup>5</sup>. Since  $q$  is a probabilistic program, scoring and sampling are simple; furthermore,  $y$  only contains the minimal number of random choices needed to specify the transition  $x \rightarrow x'$ , and thus is cheap to score. On the other hand, naive calculation of the Jacobian term and state ratio may be extremely expensive, since  $x$  may be large, for instance containing a number of variables polynomial in dataset size. As shown in algorithm 3, after sampling  $y \sim q_x$ , we factorize the move into 2 parts, both of which we provide efficient algorithms for:

1. Running  $f$  in an instrumented manner to compute the Jacobian term,  $y'$ , and an update specification  $U$  for  $x$
2. Computing  $x'$  and the state ratio from  $x$  and  $U$

### D.1. Efficiently calculating the Jacobian term

The involutive acceptance probability must include the absolute value of the determinant of the Jacobian matrix for the transformation of continuous values in  $(x, y)$  into those in  $(x', y')$ . Naively, this can be extremely expensive to calculate, since  $|x|$  may be very large, and determinant calculation is  $O(n^3)$  in matrix dimension.

**Jacobian sparsity.** In practice, most transformations  $(x, y) \rightarrow (x', y')$  only modify the values of a small number of continuous values; all other continuous values are simply copied from  $(x, y)$  to  $(x', y')$  (either at the same address or a new one). Such a copy introduces a column into the Jacobian matrix which contains all zeros except for one 1 in the position the value was copied to. Such a column cannot effect the absolute value of the matrix's determinant, so we do not need to include rows and columns corresponding to copying in the matrix. What this means in practice is that in algorithm 2, we only track transformations of the values which are explicitly read from  $x$  or  $y$  into values explicitly written into the next

---

5. In this section, we switch to using the notation  $|\det \frac{\partial(x', y')}{\partial(x, y)}|$  to mean the same thing as  $Jh_i(x, y)$  from Theorem 1, to avoid the dependency on the symbol  $h_i$ , and emphasize the idea of differentiating continuous values in  $(x', y')$  with respect to those in  $(x, y)$ .

---

**Algorithm 2** Non-standard interpretation of  $f$  for Jacobian calculation
 

---

```

procedure PROCESSINVOLUTION( $f, (x, y)$ )
    ▷ Track continuous reads and writes, manipulation moves  $M$ , and property updates  $U$ 
    ( $\text{Rd}, \text{Wr}, M, U, y'$ )  $\leftarrow \{\}, \{\}, [], \{\}, \{\}$ 
    Execute  $f$ , but with
        each manipulation command  $m$  (create, delete, change, split, merge)  $\equiv (M \leftarrow M \cup \{m\})$ 
        "proposed[ $k$ ]"  $\equiv$  (if  $y[k]$  is continuous:  $\text{Rd} \leftarrow \text{Rd} \cup \{y[k]\}; y[k]$ )
        "set backward[ $k$ ] =  $v$ "  $\equiv$  (if  $v$  is continuous:  $\text{Wr} \leftarrow \text{Wr} \cup \{y[k]\}; y' \leftarrow y' \cup \{k \mapsto v\}$ )
        "get( $P(o_1, \dots, o_n)[k]$ )"  $\equiv$  ( $v \leftarrow x[P(o_1, \dots, o_n)][k]$ ; if  $v$  is continuous:  $\text{Rd} \leftarrow \text{Rd} \cup \{v\}; v$ )
        "set  $P(o_1, \dots, o_n)[k] = v$ "  $\equiv$  (if  $v$  is continuous:  $\text{Wr} \leftarrow \text{Wr} \cup \{v\}; U[P(o_1, \dots, o_n)][k] \leftarrow v$ )
     $J \leftarrow$  uninitialized  $|\text{Rd}| \times |\text{Rd}|$  matrix
    for  $i$  in  $\{1, \dots, |\text{Wr}|\}$ :
        ▷ Execute  $f$  with reverse-mode AD to compute  $\frac{\partial \text{Wr}[i]}{\partial \text{Rd}[j]}$  for each  $j \in \{1, \dots, |\text{Rd}|\}$ 
         $J[:, i] \leftarrow \nabla_{\text{Rd}}(\text{Wr}[i])$ 
     $U \leftarrow$  CompileManipulationMoves( $U, M, x$ ) ▷ Get variable updates to implement manipulation moves.
    return ( $U, y', |\det J|$ )
end procedure
    
```

---

state  $x'$  or reversing randomness  $y'$ ; all unmentioned continuous values are simply copied without transformation. This enables us to calculate the determinant of the Jacobian matrix in  $O(1)$  in the size of  $x$  for many updates, rather than  $O(|x|^3)$ .

**Implementation via nonstandard interpretation.** Algorithm 2 shows how to use nonstandard interpretation of a world involution  $f$  to calculate the Jacobian correction term. Additionally, this interpretation tracks all object manipulation moves which are called, and all **set** commands. These are compiled into an object  $U$  which fully specifies the update  $x \rightarrow x'$ .  $U$  is a mapping from variable names  $G(o_1, \dots, o_n)$  to partial traces  $U[G(o_1, \dots, o_n)]$ , which map addresses to values.  $U$  specifies to update each trace  $x[G(o_1, \dots, o_n)]$  by overwriting addresses  $a$  in this trace with  $U[G(o_1, \dots, o_n)][a]$ . We do not detail the process of compiling object manipulation moves into update specifications; it is a fairly straightforward process of determining how much each number variable must be incremented, and changing object alias associations (which are described below).

## D.2. Efficiently calculating the state ratio

Because we use world models written using declarative **get** commands to mediate dependencies among variables in each world, our backend has access to information about which variables are conditionally independent in each state for any world model. Unlike in existing systems for programmable inference, which can only track dependency information in very simple models, this in principle enables us to efficiently calculate the state ratio

$$\frac{p(x')}{p(x)} = \frac{\prod_{v \in \text{vars}(x')} p(x'[v] | x'_{P_{a_{x'}}(v)})}{\prod_{v \in \text{vars}(x)} p(x[v] | x_{P_{a_x}(v)})} \quad (6)$$

To do this, we only need to consider variables which introduce terms to this ratio that do not cancel out. These are the variables which have been added to or deleted from  $x'$ , variables whose traces have changed between  $x$  and  $x'$ , and variables whose dependencies' values have changed.

**Complications vs updating regular Bayesian networks.** While world models are similar to Bayesian networks in that they contain a dependency graph among variables, they are far more general, and present novel challenges for producing efficient updates. In

---

**Algorithm 3** Efficient Involutive MCMC for World Models
 

---

```

procedure OUPM-IMCMC( $p, q, f, x$ )
    ( $y, q_x(y)$ )  $\leftarrow$  SAMPLEANDSCORE( $q, x$ )
    ( $U, y', J$ )  $\leftarrow$  PROCESSINVOLUTION( $f, (x, y)$ )
    ( $x', \frac{p(x')}{p(x)}$ )  $\leftarrow$  UPDATEWORLD( $x, U$ )
     $q_{x'}(y') \leftarrow$  SCORE( $q, x', y'$ )
     $\alpha \leftarrow \frac{p(x')}{p(x)} \times q_{x'}(y') \times \frac{1}{q_x(y)} \times J$ 
    with probability  $\alpha$  return  $x'$  else return  $x$ 
end procedure

 $\triangleright$  Update a trace by executing the corresponding
imperative probabilistic program.
procedure EXECUTE( $v = G(o_1, \dots, o_n), \tilde{x}, U, x$ )
     $r \leftarrow 1, C \leftarrow$  get( $U, v, \{\}$ )
     $t \leftarrow$  Trace  $G(o_1, \dots, o_n)$ , but with
    " $a \sim d$ "  $\equiv$  if  $a \in \text{vars}(C)$ :  $C[a]$ ; else:  $x[v][a]$ 
    "get  $G'(o'_1, \dots, o'_m) \equiv$  (
        if ( $u = G'(o'_1, \dots, o'_m) \notin \text{vars}(\tilde{x})$ :
             $\triangleright$  If we might have to update  $u$ , we must do
            it before accessing its value in  $v$ .
            ( $\tilde{x}, r'$ )  $\leftarrow$  EXECUTE( $G'(o'_1, \dots, o'_m), \tilde{x}, U, x$ )
             $r \leftarrow r \times r'$ 
             $\tilde{x}[u]$ 
        )
        if  $p_{\tilde{x}}(t) = 0$ :  $\triangleright$  In this case,  $v$  will be pruned.
             $\tilde{x} \leftarrow x \cup \{x[v]\}$ 
        else:
             $\tilde{x} \leftarrow \tilde{x} \cup \{v \mapsto t\}$ 
            if  $v \in \text{vars}(x)$ :
                 $r \leftarrow r \times \frac{p_{\tilde{x}}(t)}{p_x(x[v])}$ 
            else:
                 $r \leftarrow r \times p_{\tilde{x}}(t)$ 
            return ( $\tilde{x}, r$ )
        )
    )
end procedure

 $\triangleright$  Update  $x$  using  $U$  given  $\text{ord} = \text{topologicalorder}_x$ .
procedure UPDATEWORLD( $x, U, \text{ord}, \text{children}_x$ )
     $\tilde{x} \leftarrow \{\}$   $\triangleright$  New world.
     $r \leftarrow 1$   $\triangleright$  State ratio.
     $\text{pos} \leftarrow 0$   $\triangleright$  Invariant: variables topologically before
     $\text{pos}$  have been updated.
     $Q \leftarrow$  priorityqueue( $\text{ord}, \text{default} = \infty$ )
    for  $v \in \text{vars}(U)$ :
        enqueue!( $Q, v$ )
    while  $\neg \text{isEmpty}(Q)$ :
         $v \leftarrow$  pop!( $Q$ )
        if  $v \notin \text{vars}(\tilde{x})$ :
             $\triangleright$  Copy variables between  $\text{pos}$  and  $\text{ord}[v]$  to
             $\tilde{x}$ .
             $\tilde{x} \leftarrow \tilde{x} \cup \{u \mapsto x[u] : \text{ord}[u] \in [\text{pos}, \text{ord}[v]] \wedge$ 
             $u \notin \text{vars}(\tilde{x})\}$ 
             $\text{pos} \leftarrow \text{ord}[v]$ 
            ( $\tilde{x}, r'$ )  $\leftarrow$  EXECUTE( $v, \tilde{x}, U, x$ )  $\triangleright$  Update  $v$ .
             $r \leftarrow r \times r'$ 
            if  $v \in \text{vars}(x)$  and  $\text{ret}_{\tilde{x}}(\tilde{x}[v]) \neq \text{ret}_x(x[v])$ :
                for  $u \in \text{children}_x(v)$ :
                    enqueue!( $Q, u$ )
             $\triangleright$  Remove variables the observation model no
            longer depends on.
            ( $x', \text{deleted}$ )  $\leftarrow$  prune( $\tilde{x}, U$ )
            for  $v \in \text{deleted}$ :
                 $r \leftarrow r / p_{\tilde{x}}(\tilde{x}[v])$ 
                 $\tilde{x} \leftarrow \tilde{x} \upharpoonright_{\text{vars}(\tilde{x}) \setminus v}$ 
            return ( $x', r$ )
    )
end procedure
    
```

OUPMs, dependencies among variables are not fixed, and may change from state to state in a model, because the underlying Bayesian networks are contingent. While [Milch et al. \(2005b\)](#) introduced the contingent Bayesian network formalism to handle such models, they did not develop algorithms for performing arbitrary updates to these networks. Another complication is that in a world model, the values of variables in the world may deterministically depend on other variables. For instance, it would be valid for the `reading` of a seismic `Detection` to be deterministically equivalent to the `magnitude` of the corresponding `Event`. In these situations, the values of variables (`readings`) may change even if no update is specified for them, since their values are deterministic functions of other updated variables (`magnitudes`). Hence, we cannot compile an update which only visit variables in the Markov blankets of variables updates are explicitly specified for, unlike in Bayesian networks. Instead, we introduce algorithm 3, which dynamically tracks changes to the dependency graph and changes to variables which require new variables to be updated. One important insight is that when we update a variable  $a$ , we must visit every variable  $b$  which depends on  $a$  to correctly calculate the state ratio. At this time, we can check whether  $b$ 's value is deterministically changed by the update to  $a$ , and if so, enqueue the dependencies of  $b$  to be updated. This is implemented by procedure `UpdateWorld`, which pulls variables relevant to the state ratio or update from a queue, adding new variables to the

queue as needed. However, complications arise due to the necessity of updating variables in a topological order of the dependency graph.

**Update order.** The order in which variables are updated is critical for the correctness of `UpdateWorld`. For instance, say  $b \in Pa_{x'}(a)$  and  $a, b \in \text{vars}(U)$ . If we update  $b$  before we have updated  $a$ , it is impossible to calculate the factor  $p(x'[b]|x'_{Pa_{x'}(b)}) = p(x'[b]|x'[a])$  in the state ratio, because the value of  $x'[a]$  is not yet known. Thus, variables must be updated in the topological order induced by the graph of dependencies in the state  $x'$ . However, we cannot know the topological order of  $x'$  before the update, so we instead use the topological ordering of the old state  $x$ , with special-case logic to handle places where this order differs from that in  $x'$ . This is justified by the fact that efficient updates are only helpful in transitions  $x \rightarrow x'$  producing fairly local changes, and thus we expect that updates we could implement efficiently will only change the world’s topological order minorly. To implement this, the update queue  $Q$  in `UpdateWorld` is a priority queue indexed by the topological order of  $x$ . In practice, `UpdateWorld` must also update the topological order of  $x$  to one for  $x'$ ; we elide these details, along with details of maintaining the  $\text{children}_x$  lookup table (which gives the children of each variable in the CBN graph) and the details of implementing the `prune` function which removes variables which are not needed for a minimal self-supporting instantiation.

We handle changes in variable order between  $x$  and  $x'$  in the `Execute` procedure, which updates variable traces. When `Execute` observes that a variable  $v$  gets the value of a variable  $u$  which has not yet been updated, this means that in  $x$ ,  $u$  was topologically later than  $v$ , but in  $x'$ , it is topologically before  $v$ . In this situation, we recursively call `Execute(u)` to ensure that we have updated  $u$  to the proper value before it is used to update  $v$ .

**Runtime.** The full algorithm which accounts for topological order updating runs in  $O((N + M)(\log(N + M) + T))$ , where  $N$  is the number of variables which are updated or impact the state ratio,  $M$  is the number of variables whose indices in the topological ordering change during the update, and  $T$  is the maximum runtime of updating a single variable in the world.<sup>6</sup> So long as  $N \ll |\text{vars}(x')|$  and  $M \ll |\text{vars}(x')|$  (both of which are true for most local updates), this algorithm is much faster than generating a new state  $x'$  from scratch, which takes  $\theta(|\text{vars}(x')|T)$ .

**Aliased object representation.** To efficiently implement object manipulation moves, which may change the origins and indices of objects, we use a special internal representation of worlds. If we stored a mapping from the string “ $G(o_1, \dots, o_n)$ ” to traces in the world, then whenever an object  $o_k = T(\text{origin}, j)$  had its origin or index change, we would need to update the entry  $G(o_1, \dots, T(\text{origin}, j), \dots, o_n)$  for every property of this object. Instead, we assign each object a unique identifier by storing a lookup table from identifier  $i_k$  to object  $o_k = T(\text{origin}, j)$  with each world state. We then store traces in a dictionary from keys  $G(i_1, \dots, i_n)$ , where each  $i_k$  is an identifier. Index and origin changes can thus be implemented in  $O(1)$  by updating entries in the lookup table, rather than  $O(n)$  in the number of properties of the updated object.

---

6. The log terms come from the need to sort variables to update in the priority queue. Future work on static analysis of inference programs may be able to perform this sorting at compile-time in some cases.

## Appendix E. Involutive MCMC with Ancestral Regeneration.

A key benefit of involutive MCMC is that it makes it feasible to rapidly prototype inference kernels, because the acceptance-probability calculations are automated. A useful development pattern is to begin with a kernel which uses generic inference techniques, like ancestral regeneration, for proposal values for most variables, and then iterate on this kernel by adding more data-driven proposals for variables for which the generic kernels are ineffective.

To facilitate this, in the involution programming language presented in Section C, we add a command `regenerate`  $P(o_1, \dots, o_n)$  to resample property  $P$  of  $(o_1, \dots, o_n)$  from the prior distribution

$$p_{P(o_1, \dots, o_n)}(\cdot | x_{Pa_x(P(o_1, \dots, o_n))}) \quad (7)$$

where  $Pa_x(u)$  denotes the parents of variable  $u$  in the contingent Bayesian network for the open universe model the kernel operates on (with density  $p$ ), in world  $x$ .

To support this command, we must modify our formulation of automated involutive MCMC. The reason for this is that in involutive MCMC, all randomness in an inference kernel must come from the proposal  $y \sim q_x$ ;  $f$  must be a deterministic function. Therefore, we must understand the execution of `regenerate` lines in the involution program as in fact being part of the  $q_x$  proposal which produces auxiliary randomness  $y$ .

**The ancestral sampling distribution.** Given a partial world  $P$  containing values for some of the variables over which a contingent Bayesian network is defined, we define the *ancestral sampling* distribution  $r_P$  over the space  $\mathcal{W}$  of MSSIs, which completes  $P$  into a full MSSI  $w$ . If  $w$  and  $P$  disagree on any variable  $v$  (so  $w[v] \neq P[v]$ ), then  $r_P(w) = 0$ , and if they agree on all variables,

$$r_P(w) = \prod_{v \in \text{vars}(w) \setminus \text{vars}(P)} p_v(w[v] | w_{Pa_w(v)}). \quad (8)$$

**An extended proposal and involution for ancestral regeneration.** Say we have an involutive MCMC kernel defined by a partial proposal distribution  $\tilde{q}_x$  and an ‘‘involution’’ including regeneration commands  $\tilde{f}$ . To explain this as involutive MCMC, we now define a distribution  $q_x$  using  $\tilde{q}_x$  which includes the randomness from regeneration. We will describe  $q_x$  as producing tuples  $(\tilde{y}, w)$ , where  $\tilde{y} \sim \tilde{q}_x$ , and  $w$  is sampled from the ancestral resampling distribution  $r$ . Define the function  $\tilde{f}_1$  by  $\tilde{f}(x, y) = (\tilde{f}_1(x, y), \cdot)$ . We require  $\tilde{f}$  to be such that  $\tilde{f}(x, y) = \tilde{x}'$ , where  $\tilde{x}'$  is a partial world. (This comes from the user-written ‘‘involution’’  $\tilde{f}$  applying changes to world  $x$ , based on  $y$ , to produce a partial world, in which other variables may then be ancestrally sampled to produce a complete world.) We can then say  $w \sim r_{\tilde{f}_1(x, \tilde{y})}$ . We define:

$$q_x((\tilde{y}, w)) = \tilde{q}_x(\tilde{y}) \times r_{\tilde{f}_1(x, \tilde{y})}(w) \quad (9)$$

We now define an involution  $f(x, (\tilde{y}, w)) = (x', (\tilde{y}', w'))$ . For this, we require that  $\tilde{f}$  has the property that if  $\tilde{f}(x, y) = (\tilde{x}', \tilde{y}')$ , for any  $\tilde{x}' \in \mathcal{W}$  which agrees with  $\tilde{x}'$ ,  $\tilde{f}(\tilde{x}', \tilde{y}') = (\tilde{x}, \tilde{y}')$  for some  $\tilde{x}$  agreeing with  $x$ . (This is an involution in the case where  $\tilde{x}, \tilde{x}' \in \mathcal{W}$ .) Given  $f$ , we define

$$f(x, (\tilde{y}, w)) = (w, (\tilde{y}', w')) \quad (10)$$



where

$$(\tilde{x}', \tilde{y}') = \tilde{f}(x, \tilde{y}); \quad w' x_{\text{vars}(x) \setminus \text{vars}(\tilde{f}_1(\tilde{x}', \tilde{y}'))}$$

**The acceptance ratio.** Using the above distribution, for  $\tilde{f}$  satisfying similar partitioning properties as we require for involutions  $f$  in Theorem 1, the acceptance probability  $\alpha(x, (\tilde{y}, w), x', (\tilde{y}', w'))$  is

$$\frac{p(x')}{p(x)} \times \frac{\tilde{q}_{x'}(\tilde{y}') \times r_{\tilde{f}_1(x', \tilde{y}')(x)}}{\tilde{q}_x(\tilde{y}) \times r_{\tilde{f}_1(x, \tilde{y})(x')}} \times \left| \det \frac{\partial \tilde{f}(x, \tilde{y})}{\partial(x, \tilde{y})} \right| \quad (11)$$

where  $\left| \det \frac{\partial \tilde{f}(x, \tilde{y})}{\partial(x, \tilde{y})} \right|$  is the determinant of the Jacobian of the transformation of continuous values under  $\tilde{f}$  at  $(x, \tilde{y})$ .

### E.1. Automating involutive MCMC with ancestral resampling.

In section C, we introduced a language for specifying involutions on open-universe models. An involution in the language takes as arguments the current world state  $x$  and auxiliary randomness  $\tilde{y} \sim \tilde{q}_x$ , and produces a sequence  $M$  of object manipulation moves, and a set  $U$  of tuples  $(v, u)$ , where  $v$  is a variable to update the trace for, and  $u$  specifies the updated value. It also produces the reversing randomness  $\tilde{y}'$  and the Jacobian correction term  $J$ . A new state  $x'$  is constructed by applying  $M$  to the current state  $x$ , then applying all the specified updates from  $U$ . However, it is possible that the produced state  $x'$  is not a valid world model state, either because

1.  $x'$  is not minimal, in which case variables in  $x'$  are automatically deleted by our system to ensure minimality; or
2.  $x'$  is not self-supporting, in which case additional variables needed for it to be self supporting are automatically generated from the prior by our system.

Automatic deletions impose no semantic difficulties; however, when adding variables automatically, the above formulation of involutive MCMC with ancestral regeneration is needed to correctly calculate acceptance ratios.

**Accumulating the state ratio.** To implement an MCMC kernel, after sampling  $\tilde{y} \sim \tilde{q}_x$  and calculating  $(M, U, \tilde{y}', J)$ , we construct state  $x'$  using a function update  $(x, M, U)$  (Algorithm 3). To calculate the acceptance ratio, we can calculate  $\tilde{q}_{x'}(\tilde{y}')$  and  $\tilde{q}_x(\tilde{y}')$  directly, and we know  $J$ . Thus, what remains is to calculate  $\frac{p(x')}{p(x)} \times \frac{r_{\tilde{f}_1(x', \tilde{y}')(x)}}{r_{\tilde{f}_1(x, \tilde{y})(x')}}$ . When no values are generated via ancestral resampling, this reduces to the term  $\frac{p(x')}{p(x)}$ . Our implementation has the update procedure return this term as well, as it can be calculated by accumulating values while updating  $x$  to  $x'$ . We accumulate  $\frac{p(x'[v]|x'_{P_{a_{x'}}(v)})}{p(x[v]|x_{P_{a_x}(v)})}$  every time the trace for variable  $v$  is updated from  $t_v$  to  $t'_v$ ,  $\frac{1}{p(x[v]|x_{P_{a_x}(v)})}$  for every variable  $v$  which is deleted, and  $p(x'[v]|x'_{P_{a_{x'}}(v)})$  for every variable which is added. In the case where some values are produced via ancestral sampling rather than the external proposal, we must account for the  $\frac{r_{\tilde{f}_1(x', \tilde{y}')(x)}}{r_{\tilde{f}_1(x, \tilde{y})(x')}}$  term. By equation 8, the denominator of this is the product of terms  $p_v(x'[v]|x'_{P_{a_{x'}}(v)})$  for variables

$v$  which are ancestral sampled in this move from  $x$  to  $x'$ . This is straightforward for the update procedure to accumulate as it generates each of these values. The numerator of this fraction is the product of  $p_v(x[v]|x_{Pa_x(v)})$  terms for each variable  $v$  which is ancestral-sampled in the reverse move. To properly account for this, then, the update procedure needs to know which properties will be ancestral-sampled in the reversing move. We can thus define this procedure's signature as  $(x', \frac{r_{\hat{f}(x', \hat{y}')} (x)}{r_{\hat{f}(x, \hat{y})} (x')}) = \text{update}(x, M, U, R)$  where  $R$  is the set of variables which will be regenerated in the reversing move.

**Involution language extension.** To provide access to  $R$ , we extend the language from section C, adding the syntax `backward[P(o1, ..., on)] = regenerated` to specify that the reverse-direction proposal ancestral samples a value for  $P(o_1, \dots, o_n)$ . One can also write `backward[properties of (o1, ..., on)] = regenerated` to specify that all properties for any object tuple including  $o$  are regenerated by the reverse move. We also add a command `regenerate P` of  $o$ , to specify that the value for  $P(o)$  should be ancestral-sampled to produce  $x'$ , rather than copying the value for this property currently in  $x$ . (Copying from  $x$  is the default behavior when  $P(o_1, \dots, o_n)$  is present in  $x$ ; ancestral sampling is the default behavior when it is not present in  $x$ .)

One way to think about the need to specify the values ancestrally sampled in the reverse move is as follows. Although users do not specify a dictionary key for values in the full auxiliary randomness  $y$  which are generated by ancestral resampling, to ensure that the full  $f$  is a valid involution,  $f$  must still write the values of these variables to  $y'$  in the reversing move. Because users don't specify a dictionary key, we provide a special syntax, `set backward[P(o1, ..., on)] = regenerated` to automatically set the regenerated value in  $y'$  of variable  $P(o_1, \dots, o_n)$  to equal this variable's current value in  $x$ .

While the need to specify when a reverse move uses ancestral sampling adds a small cognitive burden to the user of our inference programming language, we provide dynamic error checks to ensure it is not a source of incorrectness. If users forget to include the proper `backward[...] = regenerated` statements, we can dynamically detect this error by running the reverse move and seeing which values are actually regenerated. It is thus simple to report when users have forgotten to include these commands, making it a fairly painless debugging procedure to add the proper reverse-regenerate specifications.