

## Introduction

- Rejection sampling is widely used in implementing complex generative models.
- Inference in probabilistic programs that contain unbounded loops (e.g. rejection sampling) is hard.
- We address the problem of efficient amortized importance-sampling-based inference, in particular Inference Compilation (IC) [1], in such models.
- We prove naive application of IC can produce importance weights with unbounded variance.
- We propose Amortized Rejection Sampling (ARS), an importance sampling procedure that produces unbiased expectations for programs that contain rejection sampling loops.
- We prove ARS does not introduce infinite variance in its handling of rejection sampling loops.
- We implement ARS in PyProb [2, 3] in a way that requires minimal modifications to user code.

<pre> 1: <math>x \sim p(x)</math> 2: 3: <b>for</b> <math>k \in \mathbb{N}^+</math> <b>do</b> 4:   <math>z^k \sim p(z x)</math> 5: 6: 7:   <b>if</b> <math>c(x, z^k)</math> <b>then</b> 8:     <math>z = z^k</math> 9:   <b>break</b> 10: <b>observe</b>(<math>y, p(y z, x)</math>)           (a) Original program </pre>	<pre> <math>x \sim q(x y)</math> <math>w \leftarrow \frac{p(x)}{q(x y)}</math> <b>for</b> <math>k \in \mathbb{N}^+</math> <b>do</b>   <math>z^k \sim q(z x, y)</math>   <math>w^k \leftarrow \frac{p(z^k x)}{q(z^k x, y)}</math>   <math>w \leftarrow w w^k</math>   <b>if</b> <math>c(x, z^k)</math> <b>then</b>     <math>z = z^k</math>   <b>break</b> <math>w_{IC} \leftarrow w p(y z, x)</math>           (b) Inference compilation </pre>
<pre> 1: <math>x \sim p(x)</math> 2: 3: <math>z \sim p(z x, c(x, z))</math> 4: 5: <b>observe</b>(<math>y, p(y z, x)</math>)           (c) Equivalent to above </pre>	<pre> <math>x \sim q(x y)</math> <math>w \leftarrow \frac{p(x)}{q(x y)}</math> <math>z \sim q(z x, y, c(x, z))</math> <math>w \leftarrow w \frac{p(z x, c(x, z))}{q(z x, y, c(x, z))}</math> <math>w_C \leftarrow w p(y z, x)</math>           (d) ARS </pre>

## IC weights

$$w_{IC} = \frac{p(x)}{q(x|y)} p(y|x, z) \prod_{k=1}^L \frac{p(z^k|x)}{q(z^k|x, y)}$$

**Theorem:** Under some mild conditions if the following holds then the variance of  $w_{IC}$  is infinite.

$$\mathbb{E}_{z \sim q(z|x, y)} \left[ \frac{p(z|x)^2}{q(z|x, y)^2} (1 - p(A|x, z)) \right] \geq 1$$

where  $A$  is the event of  $c(x, z)$  being satisfied.

## Collapsed weights

$$w_C = \frac{p(x)}{q(x|y)} \frac{p(z|x, A)}{q(z|x, y, A)} p(y|x, z)$$

- $\mathbb{E}[w_{IC}] = \mathbb{E}[w_C]$  but rejection sampling loops do not cause infinite variance in  $w_C$ .
- Unfortunately, we cannot directly compute  $w_C$  due to the intractable term  $\frac{p(z|x, A)}{q(z|x, y, A)}$ .

## Amortized Rejection Sampling (ARS)

$$w_C = \frac{p(x)}{q(x|y)} \frac{p(z|x)}{q(z|x, y)} p(y|x, z) \frac{q(A|x, y)}{p(A|x)}$$

- $q(A|x, y)$  is the probability of exiting the rejection sampling loop under the proposal.
- $p(A|x)$  is the probability of exiting the rejection sampling loop in the original probabilistic program.
- We use Monte Carlo sampling to get unbiased estimates of  $q(A|x, y)$  and  $\frac{1}{p(A|x)}$ .

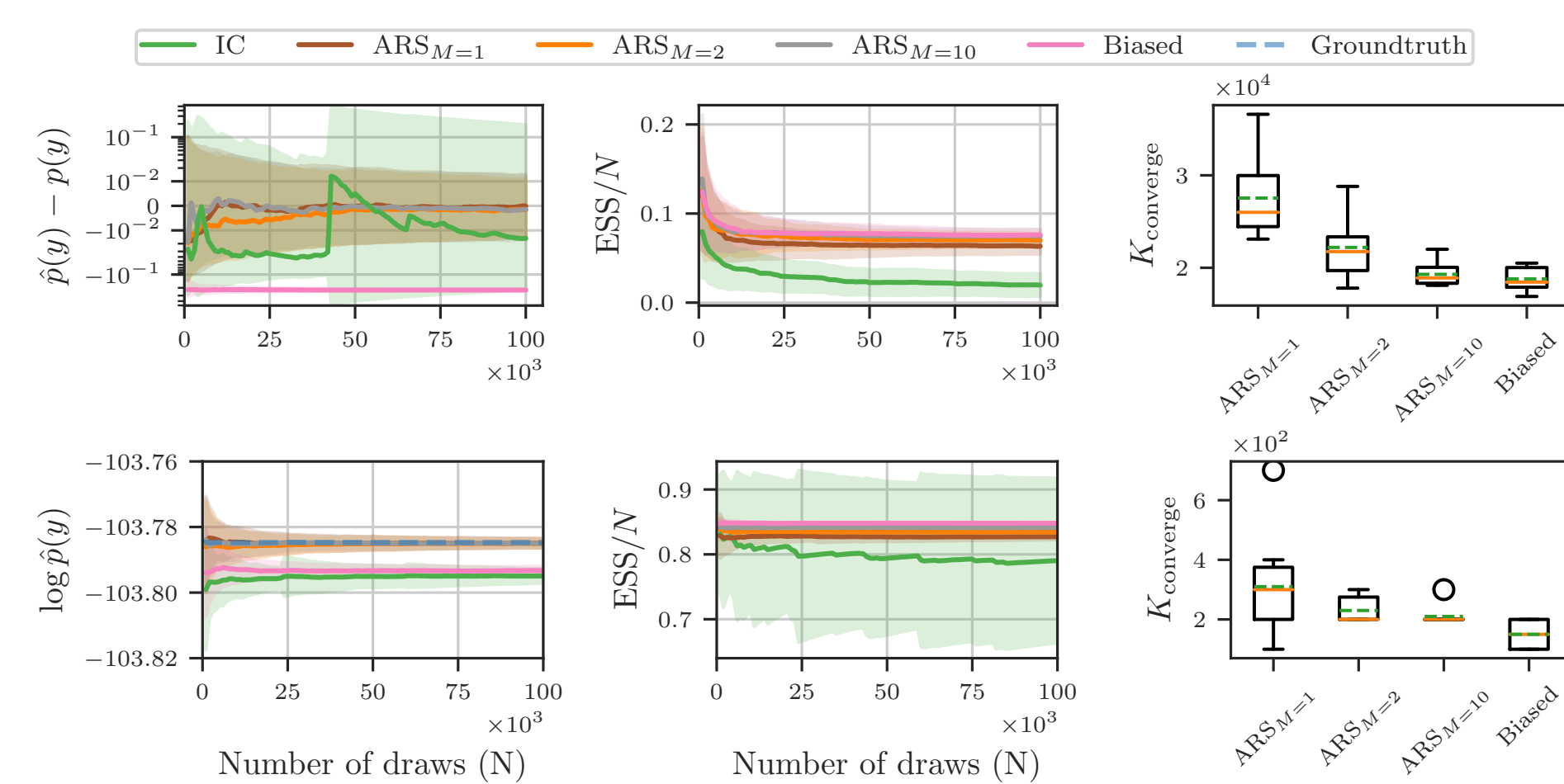
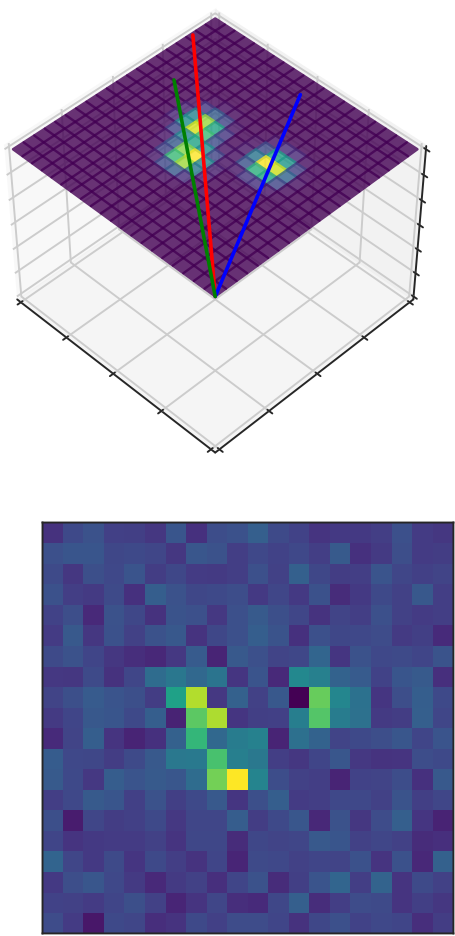
## Experiments

### Marsaglia and Mini-Sherpa

```

1: for  $k \in \mathbb{N}^+$  do
2:    $a_k \sim \text{Uniform}(-1, 1)$ 
3:    $b_k \sim \text{Uniform}(-1, 1)$ 
4:    $s = a_k^2 + b_k^2$ 
5:   if  $s < 1$  then
6:      $a = a_k$ 
7:      $b = b_k$ 
8:   break
9:  $\mu = a \sqrt{\frac{-2 \log(s)}{s}}$ 
10: observe( $y_1, \mathcal{N}(\mu, \sigma^2)$ )
11: observe( $y_2, \mathcal{N}(\mu, \sigma^2)$ )

```



(Top) Marsaglia. (Bottom) Mini-SHERPA

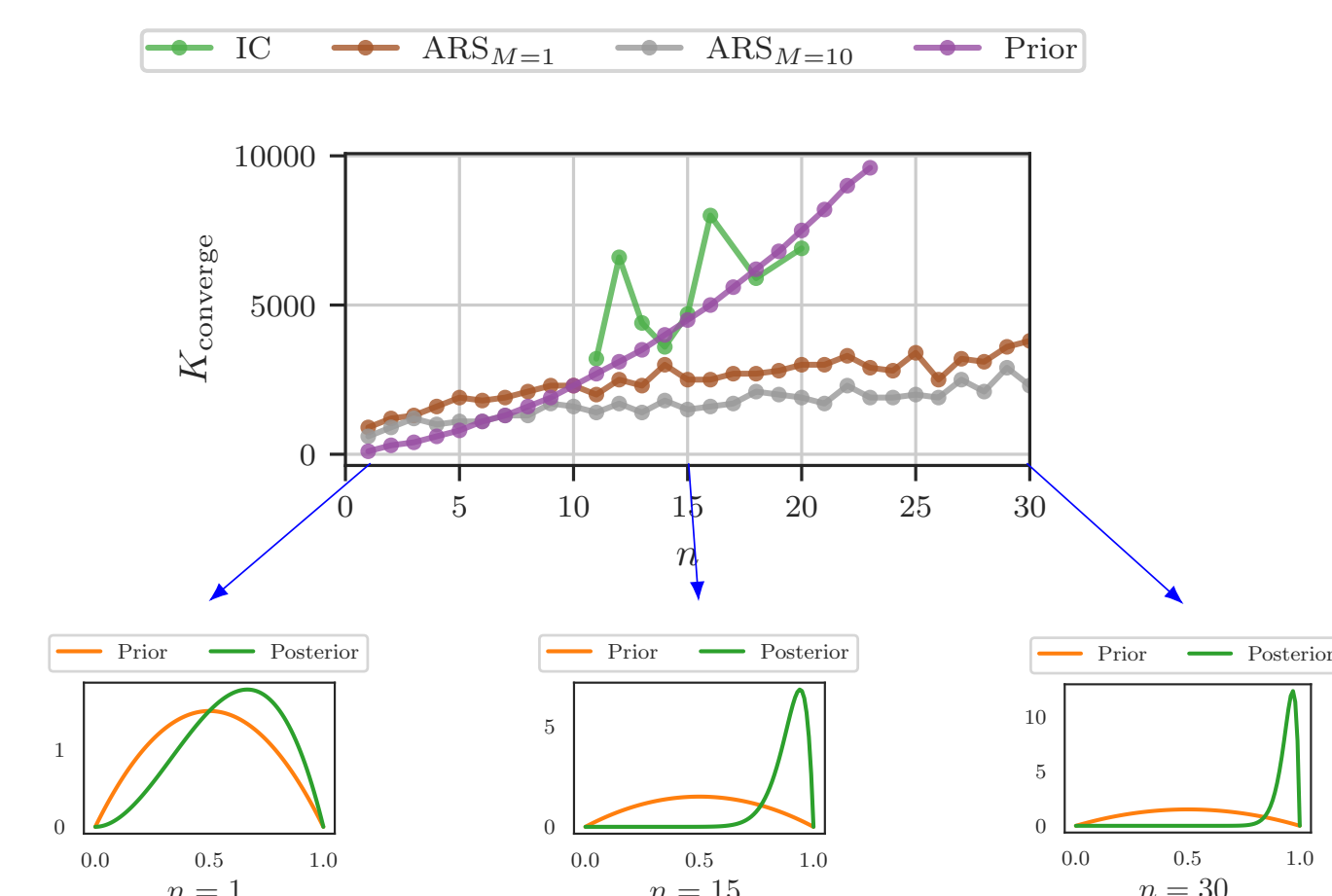
### Beta-Bernoulli

- We compare ARS with “Prior,” a baseline proposed in [2].
- This baseline uses the prior as proposal for the variables within rejection sampling loops, ignoring the learned proposal.

```

1: for  $k \in \mathbb{N}^+$  do
2:    $x_k \sim \text{Uniform}(0, 1)$ 
3:    $u_k \sim \text{Uniform}(0, 1)$ 
4:   if  $c(x_k, u_k)$  then
5:      $x = x_k$ 
6:      $u = u_k$ 
7:   break
8: for  $i \in [1, n]$  do
9:   observe(True, Ber( $x$ ))

```



- Prior converges quicker than ARS if the posterior and prior are close.
- As the difference between prior and posterior grows, ARS quickly outperforms Prior.

## Algorithm

```

1:  $x \sim q(x|y)$ 
2:  $w \leftarrow \frac{p(x)}{q(x|y)}$ 
3: for  $k \in \mathbb{N}^+$  do
4:    $z^k \sim q(z|x, y)$ 
5:   if  $c(x, z^k)$  then
6:      $z = z^k$ 
7:   break
8:  $w \leftarrow w \frac{p(z|x)}{q(z|x, y)}$ 
9:  $K \leftarrow 0$ 
10: for  $i \in 1, \dots, N$  do
11:    $z'_i \leftarrow q(z|x, y)$ 
12:    $K \leftarrow K + c(z, x)$ 
13: for  $j \in 1, \dots, M$  do
14:   for  $l \in \mathbb{N}^+$  do
15:      $z''_{j,l} \leftarrow q(z|x, y)$ 
16:     if  $c(x, z''_{j,l})$  then
17:        $T_j \leftarrow l$ 
18:     break
19:  $T \leftarrow \frac{1}{M} \sum_{j=1}^M T_j$ 
20:  $w \leftarrow w \frac{KT}{N}$ 
21:  $w \leftarrow w p(y|z, x)$ 

```

- $\frac{K}{N}$  estimates  $q(A|x, y)$
- $T$  estimates  $\frac{1}{p(A|x)}$

## Implementation

We introduce two new functions to tag the beginning and end of rejection sampling loops.

Original	Annotated
<code>x = sample(P_x)</code>	<code>x = sample(P_x)</code>
<code>while True:</code>	<code>while True:</code>
	<code>rs.start()</code>
<code>z = sample(P_z(x))</code>	<code>z = sample(P_z(x))</code>
<code>if c(x, z):</code>	<code>if c(x, z):</code>
	<code>rs.end()</code>
<code>break</code>	<code>break</code>
<code>observe(P_y(x, z), y)</code>	<code>observe(P_y(x, z), y)</code>
<code>return x, z</code>	<code>return x, z</code>

## References

- Thuan Anh Le, Atılım Güneş Baydin, and Frank Wood. Inference compilation and universal probabilistic programming. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pages 1338–1348, Fort Lauderdale, FL, USA, 2017. PMLR.
- Atılım Güneş Baydin, Lukas Heinrich, Wahid Bhimji, Bradley Gram-Hansen, Gilles Louppe, Lei Shao, Kyle Cranmer, Frank Wood, et al. Efficient probabilistic inference in the quest for physics beyond the standard model. In *Thirty-second Conference on Neural Information Processing Systems (NeurIPS)*, 2019.
- Atılım Güneş Baydin, Lei Shao, Wahid Bhimji, Lukas Heinrich, Lawrence Meadows, Jialin Liu, Andreas Munk, Saeid Naderiparizi, Bradley Gram-Hansen, Gilles Louppe, et al. Etalumis: Bringing probabilistic programming to scientific simulators at scale. In *the International Conference for High Performance Computing, Networking, Storage and Analysis (SC '19)*, 2019. doi: 10.1145/3295500.3356180.