

---

# Rejection Sampling Optimization with Parallelization, Mixture Proposals, and Whitening

---

**Aadit Rao**

aaditrao@student.ubc.ca

**Matthew Mung**

mmung3@student.ubc.ca

## Abstract

1        Rejection sampling is a simple method for sampling from difficult target distribu-  
2        tions, but its efficiency depends strongly on the proposal distribution. We study  
3        three practical ways to improve rejection sampling without changing its core frame-  
4        work: vectorized and GPU-based implementations to reduce runtime, Gaussian  
5        mixture proposals for multimodal targets, and whitening for correlated targets.  
6        Using Gaussian-based test distributions, we compare these methods against a base-  
7        line sampler using runtime, acceptance rate, and sampling efficiency. Our results  
8        show that vectorization and GPU acceleration improve throughput, while mix-  
9        ture proposals and whitening improve proposal quality and acceptance on harder  
10       targets.

## 11    1 Introduction

12    Rejection sampling is a classical method for drawing samples from a target distribution when direct  
13    sampling is difficult. Its main advantages are its simplicity and its ability to produce exact samples  
14    when a valid proposal distribution and envelope bound can be found. In practice, however, its  
15    efficiency depends on how well the proposal matches the target. When the proposal is poorly  
16    aligned with the target, many candidate samples are rejected, leading to low acceptance rates. This  
17    becomes especially problematic for multimodal and highly correlated target distributions, where  
18    simple proposals may cover large regions of low-probability space.

19    Prior work has often improved rejection sampling by making proposals more adaptive. Gilks and Wild  
20    (1992) introduced adaptive rejection sampling for univariate log-concave densities, using piecewise  
21    exponential envelopes that tighten as sampling proceeds. Gilks, Best, and Tan (1995) extended this  
22    idea to non-log-concave settings by using a Metropolis step within the rejection procedure. More  
23    recently, Raim, Livsey, and Irinata (2024) proposed vertical weighted strips, a finite-mixture proposal  
24    for weighted densities that also provides an upper bound on rejection probability. These methods  
25    show that rejection sampling can often be improved substantially through better proposal design, but  
26    they are more specialized than the standard rejection sampler considered in this project.

27    In this project, we instead study simple modifications that keep the basic rejection sampling frame-  
28    work while improving either implementation efficiency or proposal quality. We evaluate three  
29    strategies: vectorized and GPU-accelerated implementations to reduce computational overhead,  
30    Gaussian mixture proposals to better match multimodal targets, and whitening transformations to  
31    improve sampling for correlated targets. We compare these methods against a baseline rejection  
32    sampler, using runtime, acceptance rate, and sampling efficiency as evaluation metrics. Our goal is to  
33    identify which modifications are most effective under different distributional conditions.

## 34 2 Methods

### 35 2.1 Baseline Setup

36 We evaluated our methods on target distributions with increasing levels of difficulty for rejection  
37 sampling. In one dimension, we used three targets: a standard Gaussian, a bimodal Gaussian mixture,  
38 and a more complex five-component multimodal Gaussian mixture. The Gaussian target served as  
39 a simple baseline, while the bimodal and multimodal targets were chosen to highlight situations in  
40 which a single Gaussian proposal is a poor match to the target. For the whitening experiment, we  
41 used a two-dimensional correlated Gaussian target, since correlated structure is one of the main cases  
42 where rejection sampling becomes inefficient in the original coordinate system.

43 For the one-dimensional experiments, we compared three proposal strategies: a baseline single  
44 Gaussian proposal, a uniform proposal, and a Gaussian mixture proposal. For each setup, we first  
45 constructed the proposal and estimated the envelope constant  $M$  over a dense grid of 10,000 points.  
46 We then ran rejection sampling and recorded runtime, acceptance rate, and accepted sample count.  
47 Each one-dimensional experiment was repeated five times, and we reported median runtime and  
48 median acceptance statistics to reduce the effect of outliers and runtime noise. For visualization, we  
49 plotted the run whose acceptance rate was closest to the median so that the displayed results would  
50 be representative. In the standard one-dimensional experiments, each trial used 10,000 proposals.

51 For the whitening experiment, we compared direct vectorized rejection sampling in the original  
52 correlated space with sampling in a whitened space and then mapped accepted samples back to  
53 the original coordinates. This experiment used 50,000 proposals to make runtime and acceptance  
54 differences clearer. We also implemented a separate benchmarking pipeline for sequential, vectorized,  
55 and GPU-based sampling. These benchmark runs used larger sample counts of 10,000, 50,000,  
56 100,000, and 300,000 proposals, again using five trials and reporting median runtime. The CuPy-  
57 based GPU benchmarks were run on a laptop with an NVIDIA GeForce RTX 2060 GPU. This setup  
58 allowed us to evaluate proposal quality and implementation efficiency separately.

### 59 2.2 Vectorized and GPU-Accelerated Sampling

60 Our first optimization focused on improving the runtime of rejection sampling through vectorization  
61 and GPU acceleration. In the baseline implementation, candidate samples were generated and tested  
62 one at a time in a Python loop. For each proposal, the sampler drew a candidate  $x$ , generated a  
63 uniform random value  $u$ , computed the acceptance probability  $\frac{f(x)}{Mg(x)}$ , and then accepted or rejected  
64 the sample. Although this procedure is simple, it introduces substantial Python overhead when many  
65 proposals are required.

66 To reduce this cost, we implemented a vectorized NumPy version that generates a batch of proposals  
67 and random values at once, computes acceptance probabilities as arrays, and uses a boolean mask to  
68 keep accepted samples. We extended the same approach to the GPU using CuPy, replacing NumPy  
69 operations with GPU-compatible ones. This also required corresponding CuPy versions of the target  
70 density functions, since NumPy-based PDFs cannot operate directly on CuPy arrays.

71 This optimization was applied to the Gaussian, uniform, and mixture proposal samplers. Its purpose  
72 was not to change the statistical behavior of rejection sampling, but to improve computational  
73 throughput by reducing loop overhead and exploiting parallel array operations. Therefore, we  
74 expected vectorization and GPU acceleration to leave acceptance rates essentially unchanged while  
75 significantly reducing runtime, especially for larger batch sizes.

### 76 2.3 Gaussian Mixture Proposals

77 Our second optimization focused on improving proposal quality by replacing a single Gaussian  
78 proposal with a Gaussian mixture proposal. A single Gaussian can work well for simple unimodal  
79 targets, but it is often a poor fit for bimodal or multimodal distributions. In those cases, the proposal  
80 must be scaled by a larger constant  $M$  to cover all regions of the target, including the low-density  
81 space between modes. This lowers the acceptance rate and causes many samples to be wasted in  
82 regions where the target density is small.

83 To address this, we constructed the proposal as a weighted mixture of Gaussians whose components  
84 were estimated directly from the target density over a fixed range. We first evaluated the target on  
85 a dense grid and used SciPy’s peak detection routine, `scipy.signal.find_peaks`, to identify its  
86 main modes. The detected peak locations were used as the mixture means, while the peak heights  
87 were normalized to define the component weights. The component standard deviations were then  
88 estimated from the spacing between peaks and scaled slightly to provide enough overlap and coverage  
89 across the full target. If no distinct peak was found, the proposal defaulted to a single component  
90 centered at the maximum of the target density.

91 Once the mixture proposal was defined, the rest of the rejection sampling algorithm was the same.  
92 A component was first selected according to its mixture weight, then a sample was drawn from  
93 that Gaussian component, and finally the usual rejection rule was applied using the ratio between  
94 the target density and the mixture proposal density. We also computed the envelope constant  $M$   
95 numerically on a dense grid to ensure that the proposal properly bounded the target over the region of  
96 interest.

97 Unlike the vectorization and GPU optimization, this method was meant to improve acceptance rate  
98 rather than raw runtime. Its main benefit is that it shapes the proposal to better match the geometry of  
99 multimodal targets, reducing wasted sampling between modes and increasing overall efficiency. As  
100 a result, we expected this optimization to be most helpful on the bimodal and complex multimodal  
101 targets, where the limitations of a single Gaussian proposal are most apparent.

## 102 2.4 Whitening for Correlated Targets

103 Our third optimization focused on improving rejection sampling for correlated targets through  
104 whitening. Even when the target is unimodal, rejection sampling can still be inefficient if the  
105 distribution is stretched or strongly correlated across dimensions. In the original coordinate system,  
106 a simple isotropic Gaussian proposal may cover a large amount of low-probability space, which  
107 reduces acceptance rate and wastes samples.

108 To address this, we transformed the target into a whitened coordinate system where the correlation  
109 structure is removed. Using the Cholesky factor  $L$  of the covariance matrix, we mapped samples  
110 from the original space  $x$  into whitened coordinates  $y$ , where the correlated Gaussian target becomes  
111 approximately isotropic. In this transformed space, it is much easier to use a simple Gaussian  
112 proposal that aligns well with the target. After performing rejection sampling in the whitened space,  
113 the accepted samples were mapped back to the original coordinates for visualization and comparison.

114 In our implementation, we compared this whitened sampler against a baseline two-dimensional  
115 Gaussian proposal applied directly in the original correlated space. Both versions used vectorized  
116 rejection sampling, but the whitening transformation changed the geometry of the problem so that the  
117 proposal could match the target more effectively. As with the mixture proposal, this optimization  
118 was intended to improve proposal quality rather than implementation speed. Its main benefit is that it  
119 reduces inefficiency caused by correlation and anisotropy in the target distribution.

## 120 3 Experiments and Analysis

### 121 3.1 Runtime Benchmarks

122 The runtime benchmarks show that vectorization and GPU acceleration both greatly reduce the cost  
123 of rejection sampling, while leaving acceptance rates essentially unchanged. At a smaller batch size  
124 of  $n = 10,000$  (Figure 1), the vectorized NumPy implementation was the fastest across all benchmark  
125 cases. Compared with the sequential implementation, vectorization gave speedups ranging from about  
126  $65.9\times$  to  $356.5\times$ , while the GPU version still improved over sequential runtime but remained slower  
127 than NumPy vectorization. This is most likely due to GPU launch and memory overhead dominating  
128 at smaller problem sizes. For example, on the baseline Gaussian target, sequential sampling took  
129  $0.0383$  s, the vectorized version took  $5.80 \times 10^{-4}$  s, and the GPU version took  $1.20 \times 10^{-3}$  s. A  
130 similar trend appeared across the other proposal types.

131 At the larger batch size of  $n = 300,000$  (Figure 2), the results changed substantially. In this case,  
132 the GPU implementation became the fastest in every benchmark case, showing that the overhead  
133 of moving work to the GPU is outweighed once the batch size is sufficiently large. GPU speedups

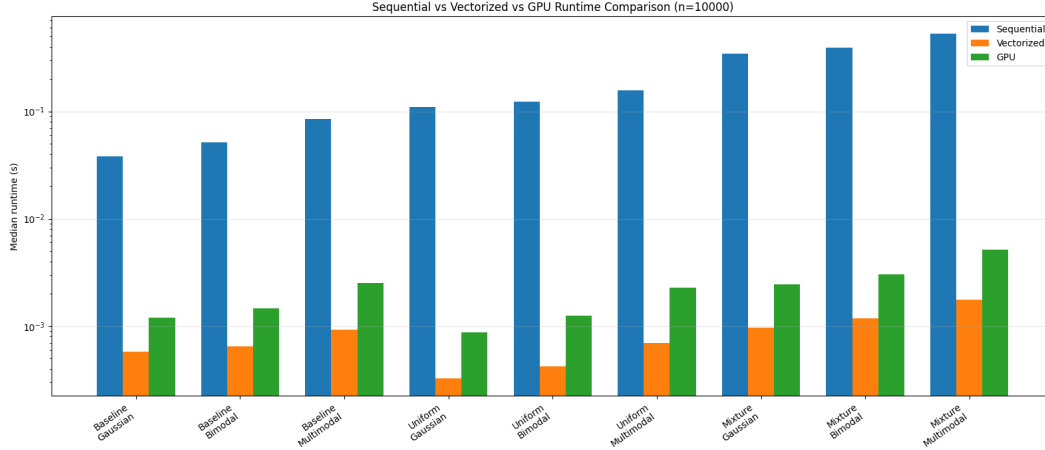


Figure 1: Runtime Benchmark n=10000

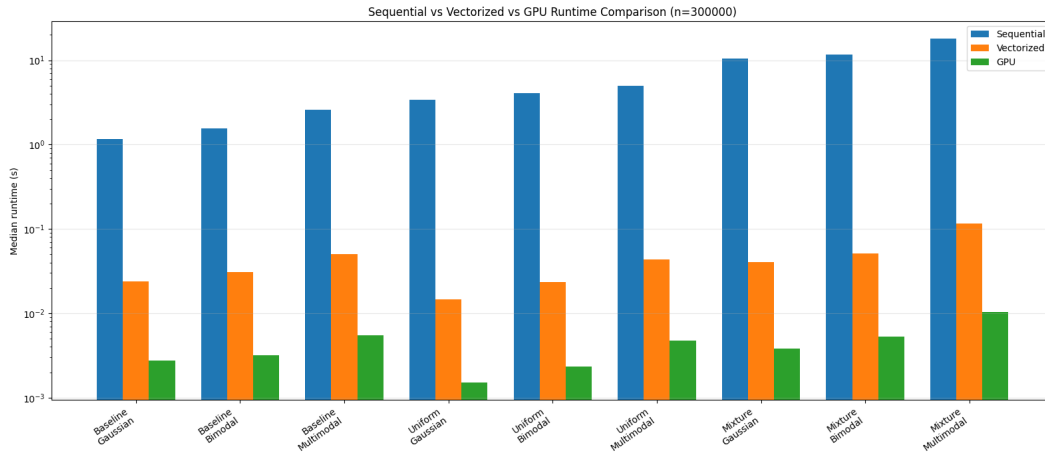


Figure 2: Runtime Benchmark n=300000

134 over the sequential implementation ranged from about  $420\times$  to  $2727\times$ , while the vectorized NumPy  
 135 version achieved speedups between about  $48\times$  and  $259\times$ . This effect was especially strong for the  
 136 more expensive uniform and mixture samplers. For instance, on the mixture multimodal benchmark,  
 137 runtime dropped from 18.01 s in the sequential case to 0.116 s with NumPy vectorization and 0.0104  
 138 s on the GPU, corresponding to a GPU speedup of about  $1731\times$ . On the uniform Gaussian benchmark,  
 139 runtime dropped from 3.43 s to 0.0146 s and then to 0.00150 s, yielding a GPU speedup of about  
 140  $2278\times$ .

141 Across all benchmark sizes, acceptance rates were nearly identical for the sequential, vectorized, and  
 142 GPU implementations. For example, at  $n = 300,000$  on the baseline Gaussian target, the acceptance  
 143 rates were 0.6351, 0.6348, and 0.6349 respectively. This consistency confirms that the optimization  
 144 changes only the implementation efficiency of the sampler rather than the underlying sampling  
 145 behavior. Overall, these results show that NumPy vectorization already provides large runtime gains  
 146 at modest batch sizes, while GPU acceleration becomes most beneficial once enough proposals are  
 147 processed to amortize device overhead.

### 148 3.2 Mixture Proposal Results

149 The mixture Gaussian proposal improved acceptance rates substantially, especially on targets with  
 150 multiple modes. While a single Gaussian proposal can work reasonably well for simple unimodal  
 151 targets, it becomes increasingly inefficient as the target develops separated peaks and irregular  
 152 structure. This trend is reflected clearly in our results. On a single Gaussian target, the acceptance

153 rate improved from 0.6338 with the single Gaussian proposal to 0.7347 with the mixture proposal.  
 154 On the bimodal target, it improved from 0.4139 to 0.5370. The largest improvement occurred on  
 155 the complex multimodal target, where the acceptance rate increased from 0.3180 to 0.5844. For  
 156 reference, the uniform proposal achieved only 0.2695 on this same target.

157 The multimodal case best illustrates why this improvement occurs. With a single Gaussian proposal  
 158 (Figure 3a), the sampler must use a large envelope constant,  $M = 3.1309$ , in order to cover all peaks  
 159 of the target distribution. As a result, many samples are drawn in low-density regions between modes  
 160 and end up being rejected. In contrast, the mixture proposal (Figure 3b) achieved a much smaller  
 161 envelope constant of  $M = 1.7050$ , showing that it matched the target density much more closely.  
 162 The detected component means matched the visible peaks of the target well, so the proposal placed  
 163 its probability mass in the right regions.

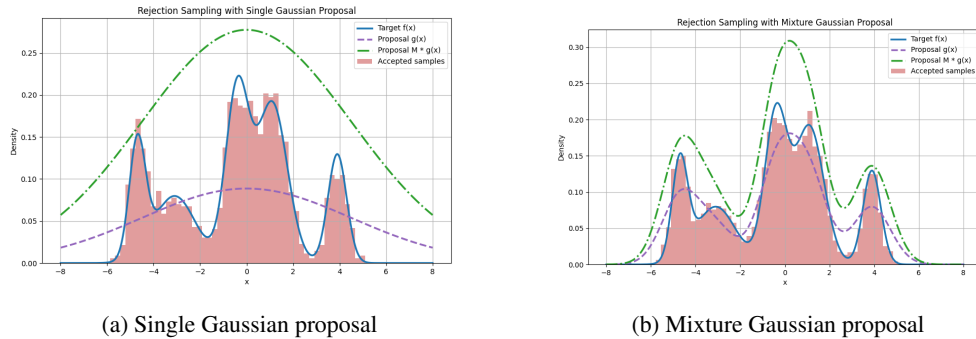


Figure 3: Rejection sampling on the complex multimodal target using a single Gaussian proposal and a Gaussian mixture proposal.

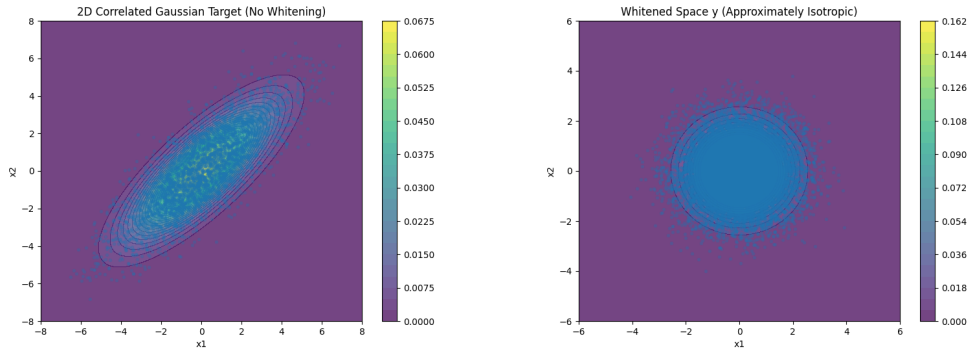
164 This difference is also visible in the sampling plots. In the single Gaussian case, the proposal is broad  
 165 and does not reflect the target's multimodal structure, so the envelope sits far above the target over  
 166 much of the domain. In the mixture case, the proposal follows the individual modes more closely, and  
 167 the histogram of accepted samples matches the target much better with less wasted space between  
 168 peaks. Although the mixture sampler was more expensive per proposal than the simpler samplers, its  
 169 main purpose was to improve the proposal itself, and the higher acceptance rates on the bimodal and  
 170 multimodal targets show that it did.

### 171 3.3 Whitening Results

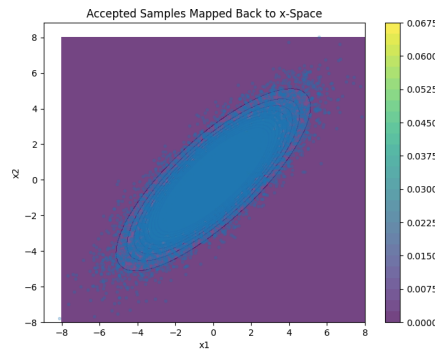
172 The whitening transformation improved rejection sampling performance substantially on the cor-  
 173 related two-dimensional Gaussian target. Without whitening, the baseline sampler in the original  
 174  $x$ -space required a large envelope constant of  $M = 10.5594$ , accepted only 4650 samples out of  
 175 50,000, and achieved an acceptance rate of 0.0930. After transforming the target into the whitened  
 176  $y$ -space, the envelope constant dropped to  $M = 2.3613$ , the number of accepted samples increased  
 177 to 21068, and the acceptance rate rose to 0.4214. Runtime also improved slightly, decreasing from  
 178 0.0152 s to 0.0106 s.

179 These results show that the main benefit of whitening was not simply computational speed, but  
 180 improved proposal quality. In the original coordinate system, the target is stretched and highly  
 181 correlated, so a simple isotropic Gaussian proposal covers a large amount of low-density space and  
 182 leads to many rejections. After whitening, the target becomes much closer to isotropic, allowing the  
 183 same style of Gaussian proposal to match the target much more closely. This is reflected in the much  
 184 smaller value of  $M$  and the large increase in acceptance rate.

185 The geometric effect of whitening is also clear in the plots. In the original  $x$ -space (Figure 4a), the  
 186 target has a long tilted elliptical shape, which makes it difficult for an isotropic proposal to cover  
 187 efficiently. In the whitened space (Figure 4b), the target becomes approximately circular, making  
 188 rejection sampling much more efficient. After the accepted samples are mapped back to the original  
 189 coordinates (Figure 4c), they again match the correlated elliptical structure of the original target. This  
 190 confirms that whitening improves efficiency during sampling while still producing samples in the  
 191 correct original distribution.



(a) Direct rejection sampling in the original  $x$ -space without whitening. (b) The same target represented in the whitened  $y$ -space.



(c) Accepted samples from the whitened sampler mapped back to the original  $x$ -space.

Figure 4: Effect of whitening on rejection sampling for a correlated two-dimensional Gaussian target

## 192 4 Discussion and Future Work

### 193 4.1 Main Conclusions

194 In this project, we studied several practical ways to improve the efficiency of rejection sampling  
 195 without changing its underlying framework. Starting from a baseline rejection sampler, we evaluated  
 196 three optimizations: vectorized and GPU-accelerated sampling, Gaussian mixture proposals, and  
 197 whitening for correlated targets. Each method addressed a different weakness of standard rejection  
 198 sampling, allowing us to study both implementation efficiency and proposal quality.

199 Our results showed that vectorization and GPU acceleration greatly reduced runtime while leaving  
 200 acceptance rates essentially unchanged. The mixture Gaussian proposal improved acceptance rates  
 201 substantially on bimodal and multimodal targets by matching the target structure more closely than a  
 202 single Gaussian proposal. Whitening also proved effective for correlated two-dimensional targets,  
 203 where transforming the distribution into a less correlated space led to a much smaller envelope  
 204 constant and a large increase in acceptance rate.

205 Overall, these results show that rejection sampling can be improved significantly through relatively  
 206 simple but well-motivated modifications. Rather than proposing a new sampling method, our work  
 207 demonstrates that better proposal design and more efficient implementations can make classical  
 208 rejection sampling much more practical in difficult settings.

## 209 4.2 Strengths and Limitations

210 A strength of this project is that it compares several distinct types of improvements within a common  
211 experimental framework. Instead of focusing only on faster implementation or better proposal design,  
212 we studied both and showed how they affect different aspects of rejection sampling performance.  
213 This made it possible to separate improvements in computational throughput from improvements in  
214 acceptance rates.

215 One limitation of this project is that our experiments were carried out on a relatively small set of  
216 target distributions. While these targets were chosen to highlight important weaknesses of standard  
217 rejection sampling, they were still mostly low-dimensional and analytically simple. In particular,  
218 our whitening experiment focused only on a two-dimensional correlated Gaussian, and our mixture  
219 proposal experiments were limited to one-dimensional multimodal targets.

220 Another limitation is that the performance study was conducted on a limited range of hardware and  
221 benchmark settings. Our GPU experiments were run on a computer with an NVIDIA GeForce RTX  
222 2060 GPU, so the runtime results may differ on stronger desktop GPUs or on systems with different  
223 memory and transfer overheads. We also used five trials per setting and reported median runtime and  
224 acceptance statistics, which gave stable comparisons, but more trials would allow for more thorough  
225 results on runtime variability.

226 A further limitation is that our proposal construction methods were still fairly simple. For the mixture  
227 proposal, the Gaussian components were estimated using peak detection and simple heuristics for the  
228 component widths. For whitening, we assumed that the covariance structure of the target was known  
229 or could be specified directly. These choices were reasonable, but they limit how broadly the results  
230 can be generalized to more realistic sampling problems.

## 231 4.3 Future Directions

232 With more time, we would test these methods on a wider range of target distributions, including  
233 higher-dimensional correlated targets, more irregular multimodal distributions, and targets whose  
234 structure is less cleanly captured by Gaussian components. This would help determine how well the  
235 observed improvements behave beyond the relatively simple distributions considered here.

236 We would also expand the benchmarking across different CPUs and GPUs, test larger sample sizes,  
237 and measure additional performance metrics such as accepted samples per second and the cost of  
238 host-device data transfer separately. Together, these extensions would provide a clearer picture of  
239 when the proposed optimizations are most beneficial in practice.

## 240 References

- 241 W. R. Gilks, N. G. Best, and K. K. C. Tan (1995). “Adaptive Rejection Metropolis Sampling Within  
242 Gibbs Sampling.” *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 44.4,  
243 pp. 455–472. DOI: 10.2307/2986138.
- 244 W. R. Gilks and P. Wild (1992). “Adaptive Rejection Sampling for Gibbs Sampling.” *Journal of the*  
245 *Royal Statistical Society: Series C (Applied Statistics)* 41.2, pp. 337–348. DOI: 10.2307/2347565.
- 246 Andrew M. Raim, James A. Livsey, and Kyle M. Irimata (2024). “Rejection Sampling with Vertical  
247 Weighted Strips.” DOI: 10.48550/arXiv.2401.09696. arXiv: 2401.09696.