

# The premise of approximate MCMC in Bayesian deep learning

Theodore Papamarkou

N8CIR

Leeds

1 November 2022

# Talk structure

- 1 Problem specification
- 2 Possible ways forward
- 3 Future work

## Topic and motivation

- Topic: approximate MCMC methods for neural networks
- Motivation: uncertainty quantification for neural networks

# Status quo of Bayesian inference in high dimensions

	<b>ML</b>	<b>Statistics</b>
Approaches	Approximate models Ex: variational inference Experimental validation	Transformed models Ex: stereographic MCMC Exact inference (convergence)
Suggestions	Decompose exact model	Predictive performance bounds

## Understanding what is the problem before solving it

- ML/statistics attack the problem computationally/theoretically
- But what is the problem really?
- **Zero acceptance rates** are observed
- Theory: lack of convergence, symmetries, identifiability
- Practically, what are the roots of the problem?

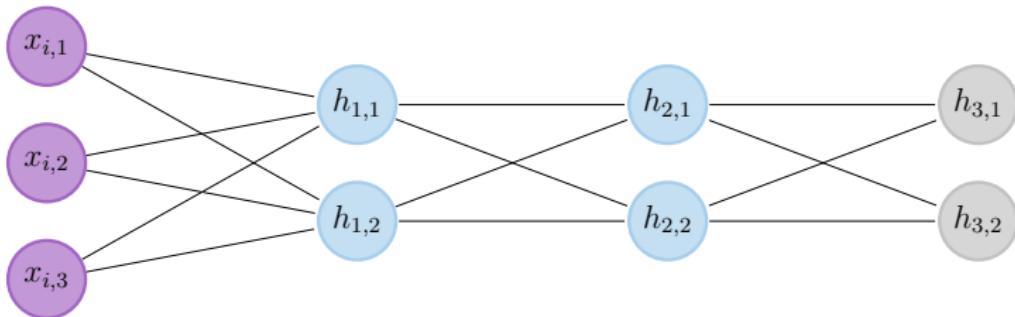
## Issue 1: big data (number of samples)

Dataset		
Name	Sample size	
	Training	Test
MNIST	60,000	10,000
Fashion MNIST (FMNIST)	60,000	10,000

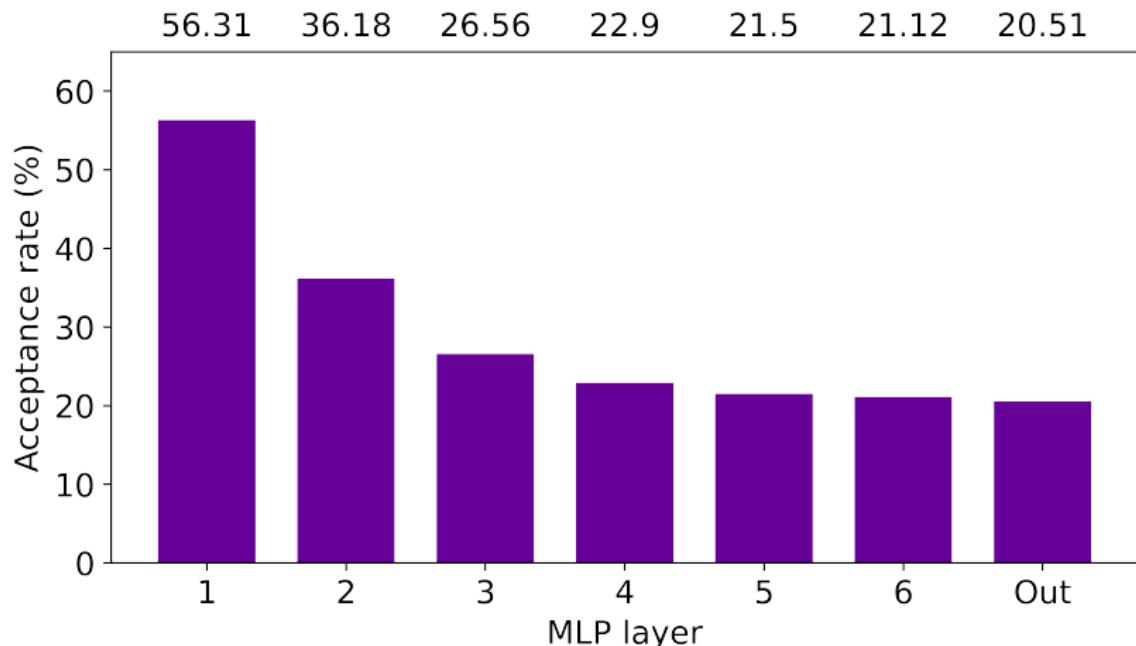
## Issue 2: big models (number of parameters)

Neural network	
Architecture	# parameters
MLP(784, 10, 10, 10, 10)	8,180

## Issue 3: hierarchical model structure (part I)



## Issue 3: hierarchical model structure (part II)



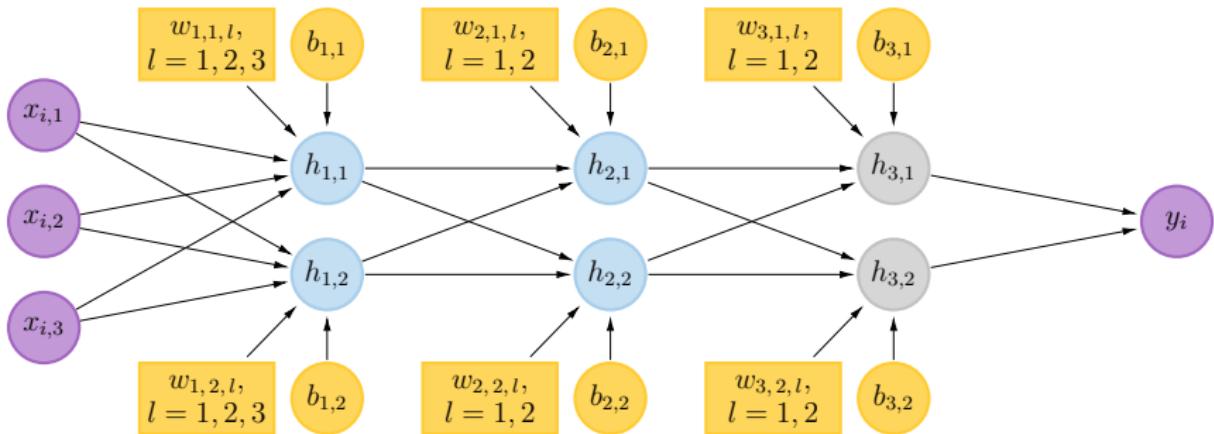
# Summary of issues

Issue	
1	Number of samples
2	Number of parameters
3	Hierarchical model

# Possible solutions

	<b>Issue</b>	<b>Solution</b>
1	Number of samples	Minibatch sampling
2	Number of parameters	Block sampling
3	Hierarchical model	Depth ↗ ⇒ Step ↘

# Training (part I)



# Training (part II)

---

## Algorithm 1 Metropolis-within-blocked-Gibbs sampling

---

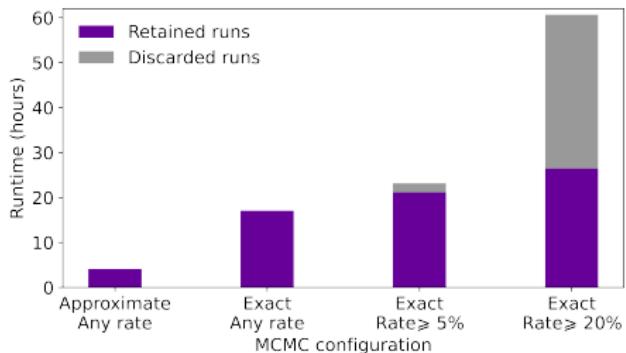
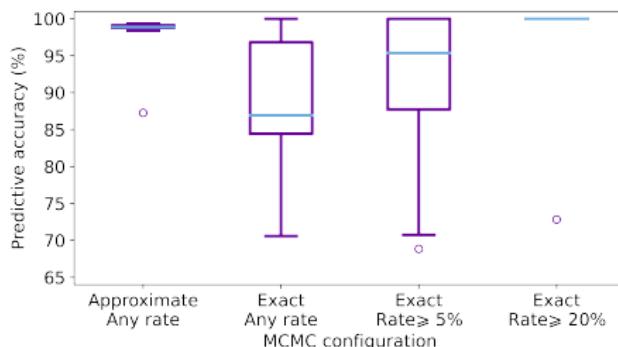
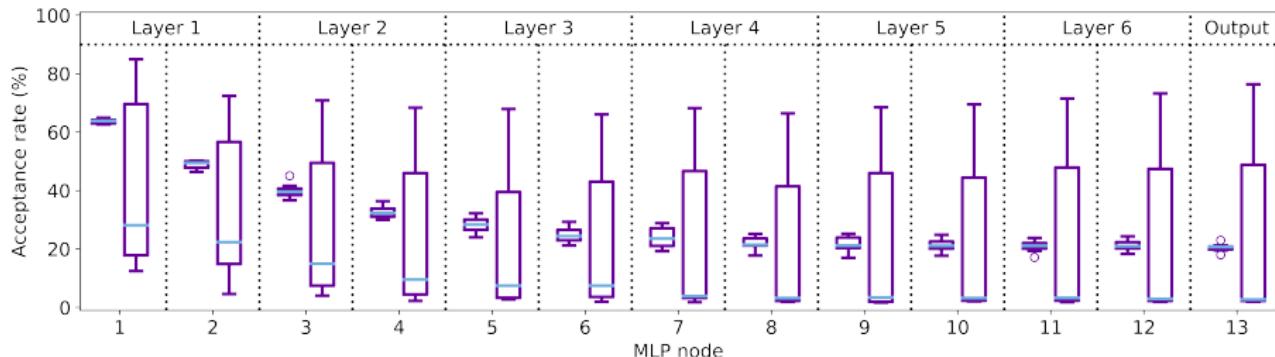
```
1: Input: training dataset  $\mathcal{D}_{1:s}$ 
2: Input: initial state  $\theta_{z(1):z(m)}^{(0)}$ 
3: Input: proposal variances  $(\sigma_1^2, \dots, \sigma_m^2)$  across blocks

4: for  $t = 1, \dots, v$  do
5:   for  $q = 1, \dots, m$  do
6:      $\theta_{z(q)}^* \sim \mathcal{N}(\theta_{z(q)}^{(t-1)}, \sigma_q^2 I_q)$ 
7:      $a(\theta_{z(q)}^*, \theta_{z(q)}^{(t-1)}) = \min \left\{ \frac{\pi(\theta_{z(q)}^*)}{\pi(\theta_{z(q)}^{(t-1)})} \exp \left( \mathcal{E}(\theta^{(t-1)}, \mathcal{D}_{1:s}) - \mathcal{E}(\theta^*, \mathcal{D}_{1:s}) \right), 1 \right\}$ 
8:      $u \sim \mathcal{U}(0, 1)$ 
9:     if  $u \leq a(\theta_{z(q)}^*, \theta_{z(q)}^{(t-1)})$  then
10:       Set  $\theta_{z(q)}^{(t)} = \theta_{z(q)}^*$ 
11:     else
12:       Set  $\theta_{z(q)}^{(t)} = \theta_{z(q)}^{(t-1)}$ 
13:     end if
14:   end for
15: end for
```

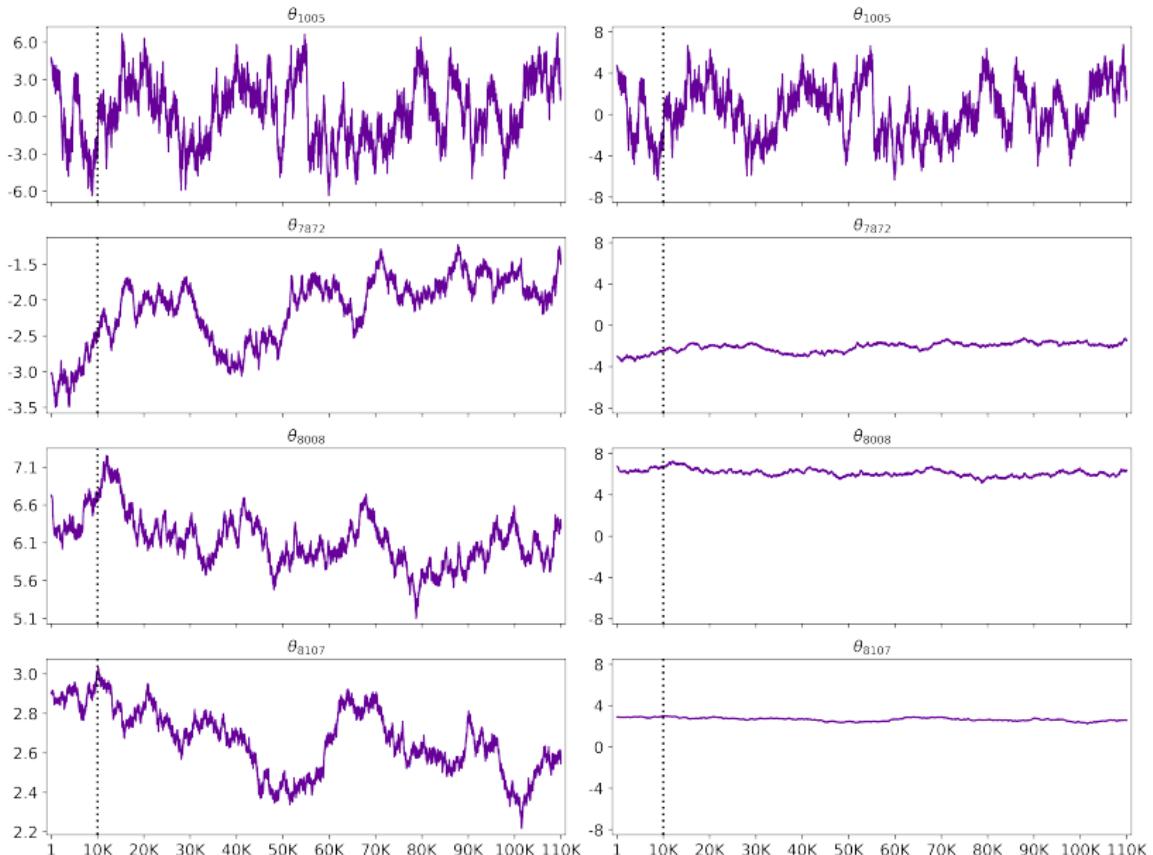
## Testing (Bayesian marginalization)

$$\underbrace{p(y|x, D_{1:s})}_{\text{Predictive distribution}} = \int \underbrace{p(y|x, \theta)}_{\text{Likelihood}} \underbrace{p(\theta|D_{1:s})}_{\text{Parameter posterior}} d\theta$$

# Revisiting issue 1: exact vs minibatch MCMC



## Revisiting issue 2: finer node-blocked sampling (part I)



## Revisiting issue 2: finer node-blocked sampling (part II)

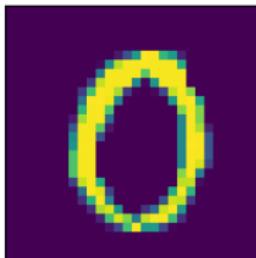
Dataset	Chain length			
	1K	10K	20K	30K
MNIST	88.31	90.75	91.12	91.20
FMNIST	78.93	80.89	81.36	81.53

## Revisiting issue 3: acceptance rates and depth (MNIST)

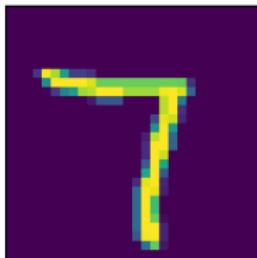
	Layer	$\sigma$	Rate
Batch size = 1800 (3%)			
Hidden	1 <sup>st</sup>	$2 \cdot 10^{-2}$	41.41
	2 <sup>nd</sup>	$2 \cdot 10^{-4}$	30.68
	3 <sup>rd</sup>	$2 \cdot 10^{-4}$	31.92
	Output	$2 \cdot 10^{-5}$	35.66
Batch size = 3000 (5%)			
Hidden	1 <sup>st</sup>	$10^{-2}$	54.95
	2 <sup>nd</sup>	$10^{-4}$	45.73
	3 <sup>rd</sup>	$10^{-4}$	44.98
	Output	$10^{-5}$	51.54

# Uncertainty quantification

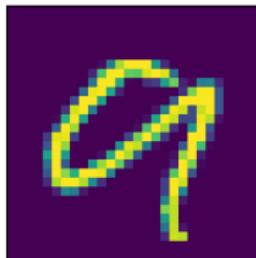
0: 0.98



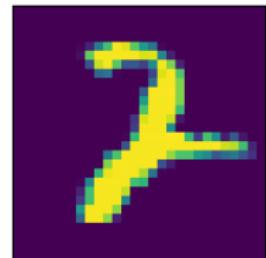
7: 0.97



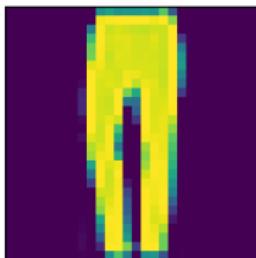
9: 0.35  
4: 0.28



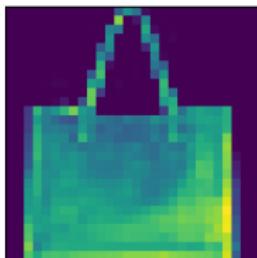
2: 0.58  
1: 0.25



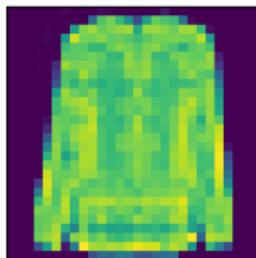
Trousers: 0.99



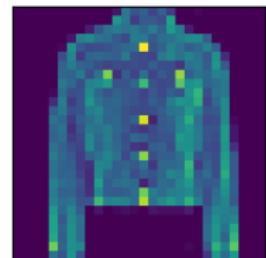
Bag: 0.96



Shirt: 0.33  
Pullover: 0.32



Coat: 0.37  
Shirt: 0.34



## Future work

- Automate tuning of proposal variances via adaptive Gibbs
- Examine alternative parameter blocks
- Derive lower bounds of predictive accuracy