

Understanding **weight matrices** and **parameters count** of different deep learning ensemble models used for **Uncertainty Estimation** –  
**Deep Ensemble, Batch Ensemble & Rank-1 BNN**

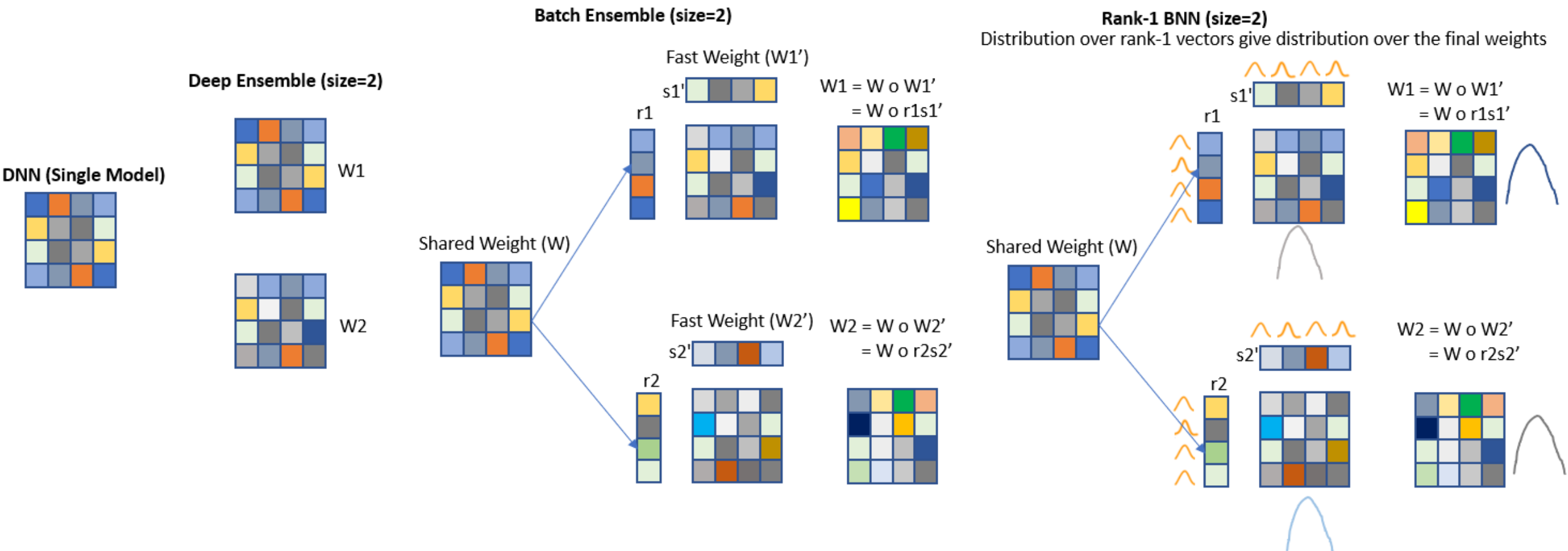
Anuj Shah

## **Part-1**

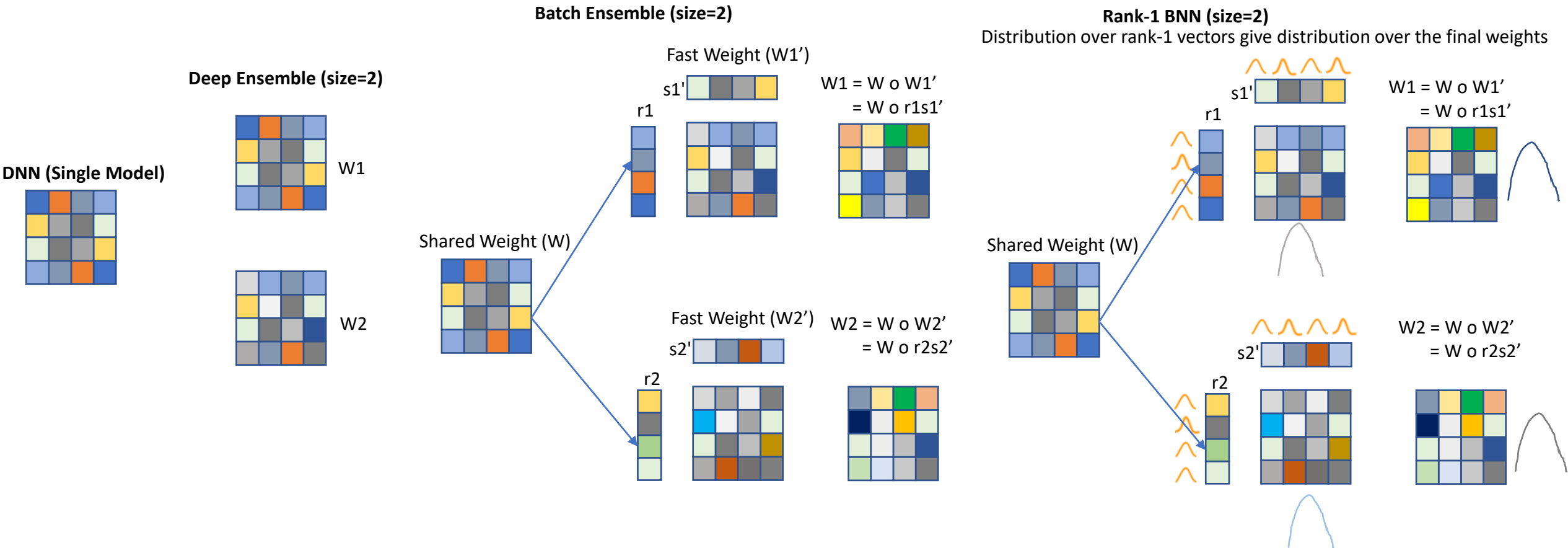
# Understanding **weight matrices** and **parameters count** of different deep learning ensemble models used for **Uncertainty Estimation**

## Part-1: Theoretical understanding

### Deep Ensemble vs Batch Ensemble vs Rank-1 BNN



# Deep Ensemble vs Batch Ensemble vs Rank-1 BNN



---

## Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles

---

Balaji Lakshminarayanan   Alexander Pritzel   Charles Blundell  
DeepMind  
{balaji, ap, cblundell}@google.com

### Abstract

Deep neural networks (NNs) are powerful black box predictors that have recently achieved impressive performance on a wide spectrum of tasks. Quantifying predictive uncertainty in NNs is a challenging and yet unsolved problem. Bayesian NNs, which learn a distribution over weights, are currently the state-of-the-art for estimating predictive uncertainty; however these require significant modifications to the training procedure and are computationally expensive compared to standard (non-Bayesian) NNs. We propose an alternative to Bayesian NNs that is simple to implement, readily parallelizable, requires very little hyperparameter

## BATCHENSEMBLE: AN ALTERNATIVE APPROACH TO EFFICIENT ENSEMBLE AND LIFELONG LEARNING

Yeming Wen<sup>1,2,3\*</sup>, Dustin Tran<sup>3</sup> & Jimmy Ba<sup>1,2</sup>

<sup>1</sup>University of Toronto, <sup>2</sup>Vector Institute, <sup>3</sup>Google Brain

### ABSTRACT

Ensembles, where multiple neural networks are trained individually and their predictions are averaged, have been shown to be widely successful for improving both the accuracy and predictive uncertainty of single neural networks. However, an ensemble's cost for both training and testing increases linearly with the number of networks, which quickly becomes untenable.

In this paper, we propose BatchEnsemble<sup>1</sup>, an ensemble method whose computational and memory costs are significantly lower than typical ensembles. BatchEnsemble achieves this by defining each weight matrix to be the Hadamard product of a shared weight among all ensemble members and a rank-one matrix per member. Unlike ensembles, BatchEnsemble is not only parallelizable across devices, where one device trains one member, but also parallelizable within a device, where multiple ensemble members are updated simultaneously for a given mini-batch. Across CIFAR-10, CIFAR-100, WMT14 EN-DE/EN-FR translation, and out-of-distribution tasks, BatchEnsemble yields competitive accuracy and uncertainties as typical ensembles; the speedup at test time is 3X and memory reduction is 3X at an ensemble of size 4. We also apply BatchEnsemble to lifelong

---

## Efficient and Scalable Bayesian Neural Nets with Rank-1 Factors

---

Michael W. Dusenberry<sup>\*†1</sup> Ghassen Jerfel<sup>\*12</sup> Yeming Wen<sup>13</sup> Yi-An Ma<sup>14</sup> Jasper Snoek<sup>1</sup>  
Katherine Heller<sup>12</sup> Balaji Lakshminarayanan<sup>1</sup> Dustin Tran<sup>1</sup>

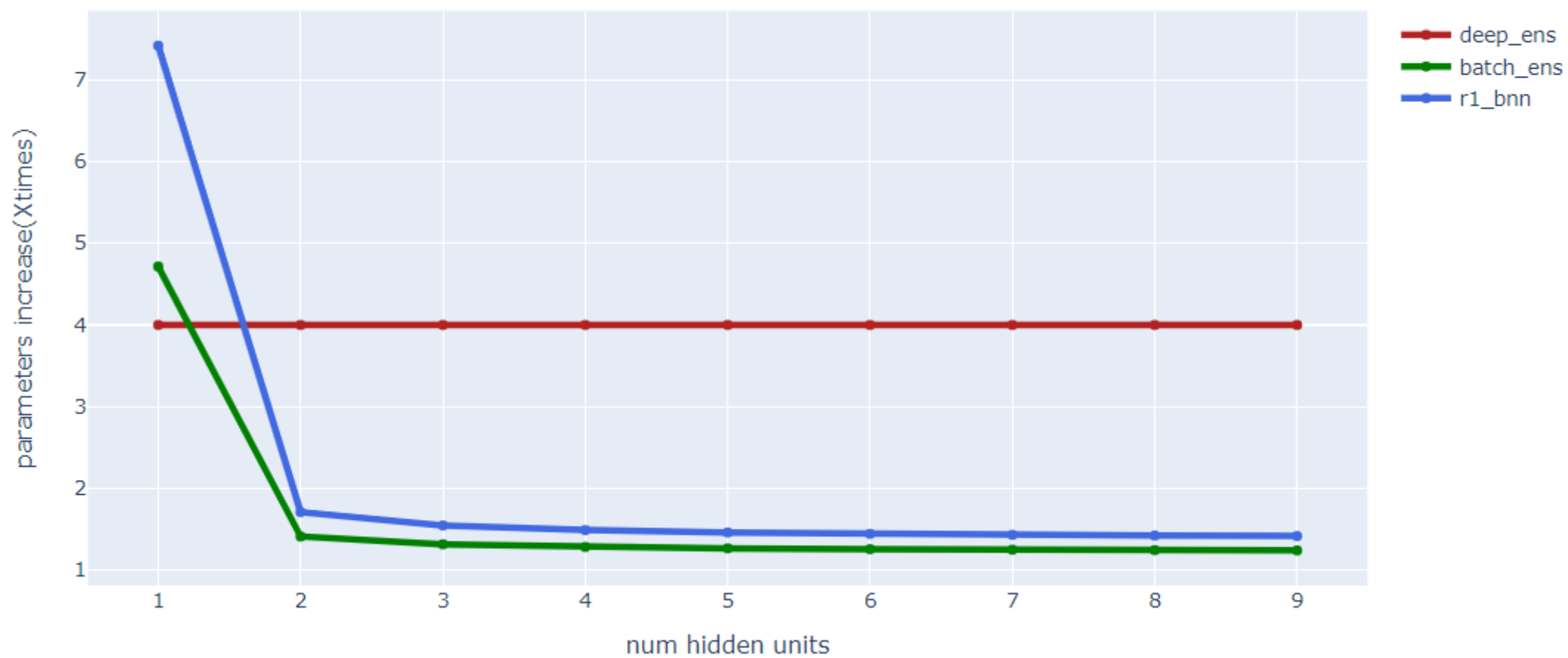
### Abstract

Bayesian neural networks (BNNs) demonstrate promising success in improving the robustness and uncertainty quantification of modern deep learning. However, they generally struggle with underfitting at scale and parameter efficiency. On the other hand, deep ensembles have emerged as alternatives for uncertainty quantification that, while outperforming BNNs on certain problems, also suffer from efficiency issues. It remains unclear how to combine the strengths of these two

deep learning. In principle, BNNs can permit graceful failure, signalling when a model does not know what to predict (Kendall & Gal, 2017; Dusenberry et al., 2019), and can also generalize better to out-of-distribution examples (Louizos & Welling, 2017; Malinin & Gales, 2018). However, there are two important challenges prohibiting their use in practice.

First, Bayesian neural networks often underperform on metrics such as accuracy and do not scale as well as simpler baselines (Gal & Ghahramani, 2016; Lakshminarayanan et al., 2017; Maddox et al., 2019). A possible reason is that the best configurations for BNNs remain unknown. What

param\_comparison\_50 (ensemble: 4)



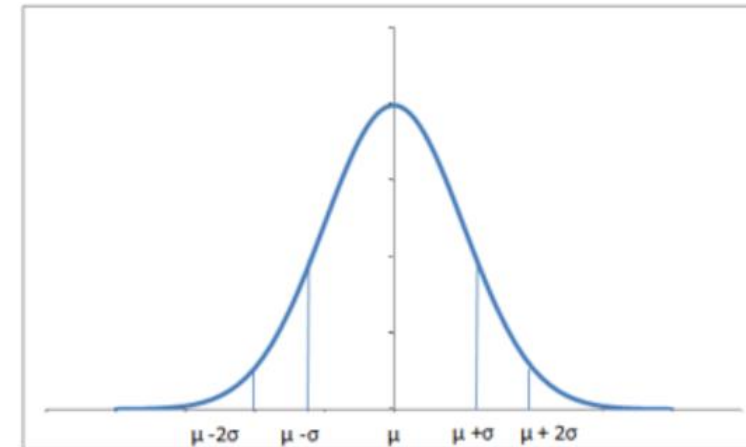
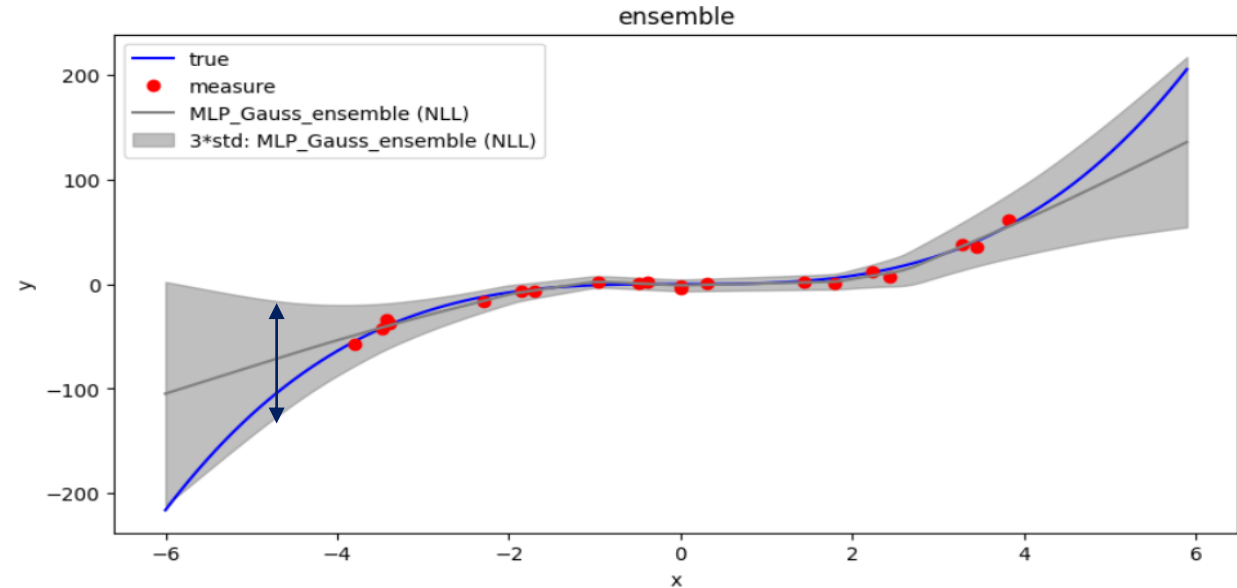
# Uncertainty Estimation

- In this video we are not going to cover in depth about the uncertainty estimation or the ensemble approaches
- We will see for the three approaches
  - On a high level how the weight matrices are defined
  - Code a simple MLP for deep ensemble , batch ensemble and rank-1 BNN
  - Compute the parameters for all the above
  - then convince us mathematically about the number of parameters in each of this paper make sense
  - and show the comparison through graph

# Uncertainty Estimation

- Uncertainty estimation refers to quantifying the reliability of your prediction.
- A result can be uncertain even when it is highly confident

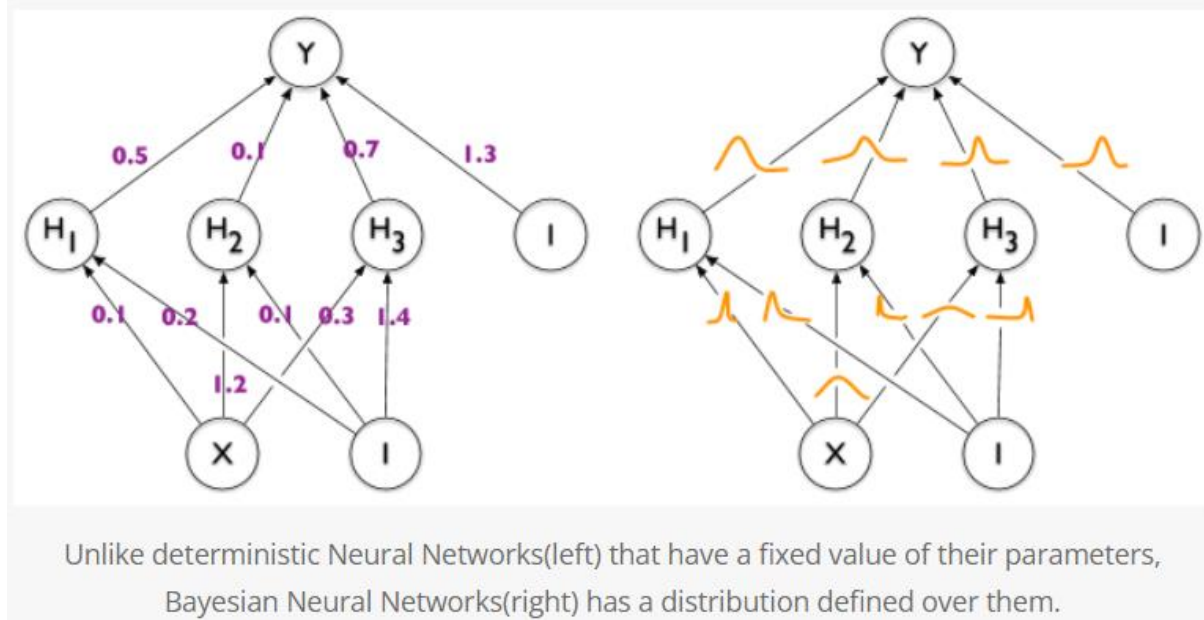
**Prediction uncertainty** refers to the variability in **prediction** due to plausible alternative input values.





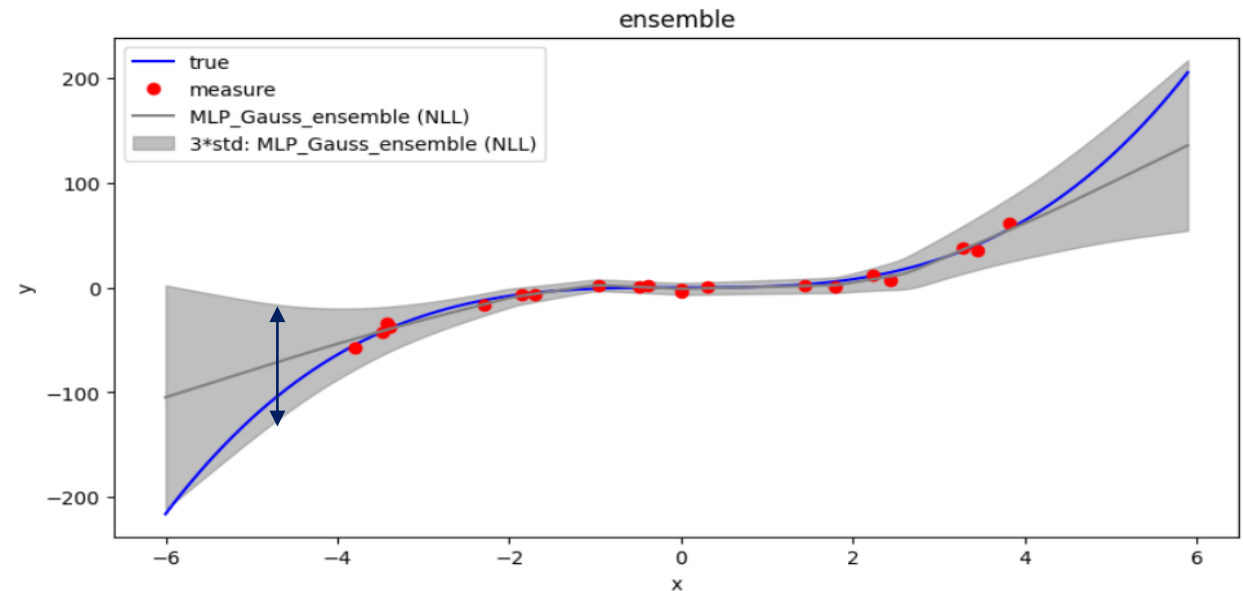
# Uncertainty Estimation

- Bayesian Neural networks can estimate predictive uncertainty but in practice BNN are harder to implement and train
- They are also computationally expensive



# Uncertainty Estimation

- The Paper *“Simple & Scalable Predictive Uncertainty Estimation using Deep Ensembles – L Balaji”*, showed that ensemble in deep learning is a good way of uncertainty estimation and can serve as an alternative to the Bayesian neural networks.
- Various papers followed this and proposes approaches for efficient ensemble. Two of which are – *“BatchEnsemble: An Alternative Approach to Efficient Ensemble and Lifelong Learning – Yeming Wen”* & *“Efficient and Scalable Bayesian Neural Nets with Rank-1 Factors – M.W. Dusenberry & G Jerfel”*

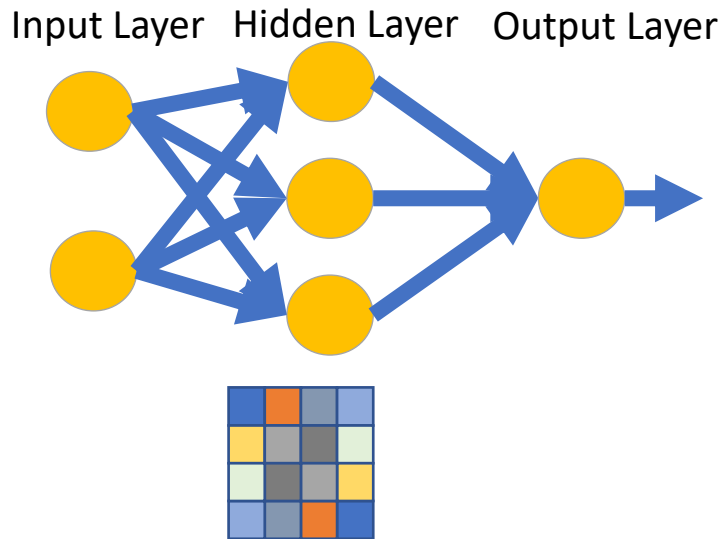


# Preface

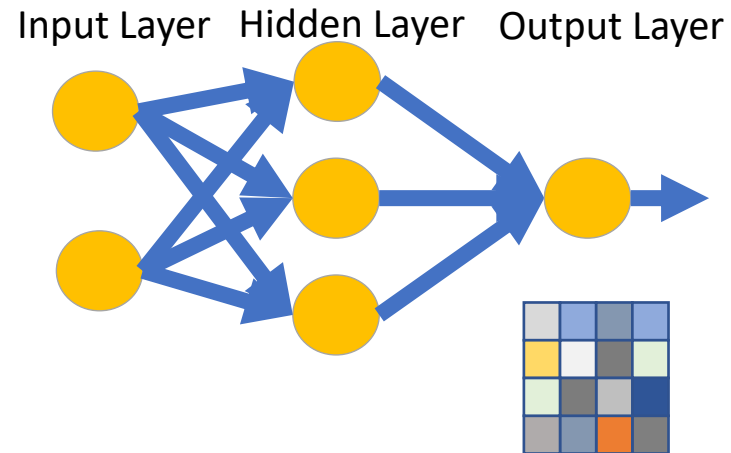
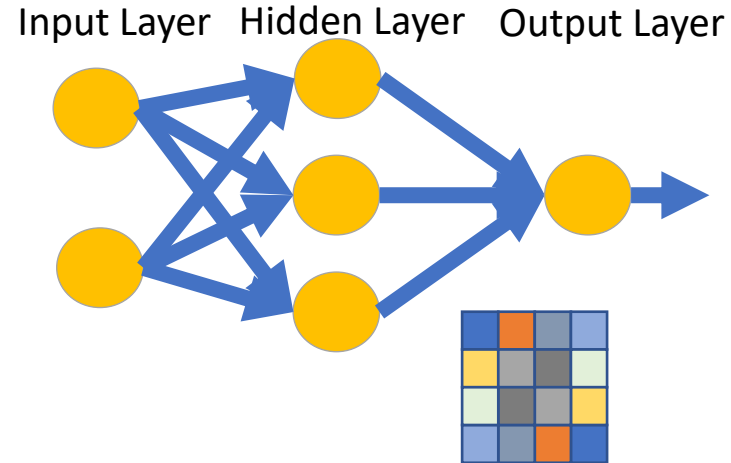
- Simple & Scalable Predictive Uncertainty Estimation using Deep Ensembles – L Balaji : **Deep Ensemble**.
- BatchEnsemble: An Alternative Approach to Efficient Ensemble and Lifelong Learning – Yeming Wen : **Batch Ensemble**.
- Efficient and Scalable Bayesian Neural Nets with Rank-1 Factors – M.W. Dusenberry & G Jerfel : **Rank-1 BNN**.
- Ensemble Size in explanation -2
- Ensemble Size in code and parameter count -4

# Deep Ensemble

DNN (Single Model)

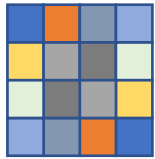


Deep Ensemble (size=2)

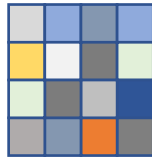
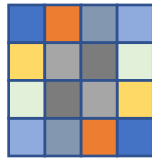


# Deep Ensemble

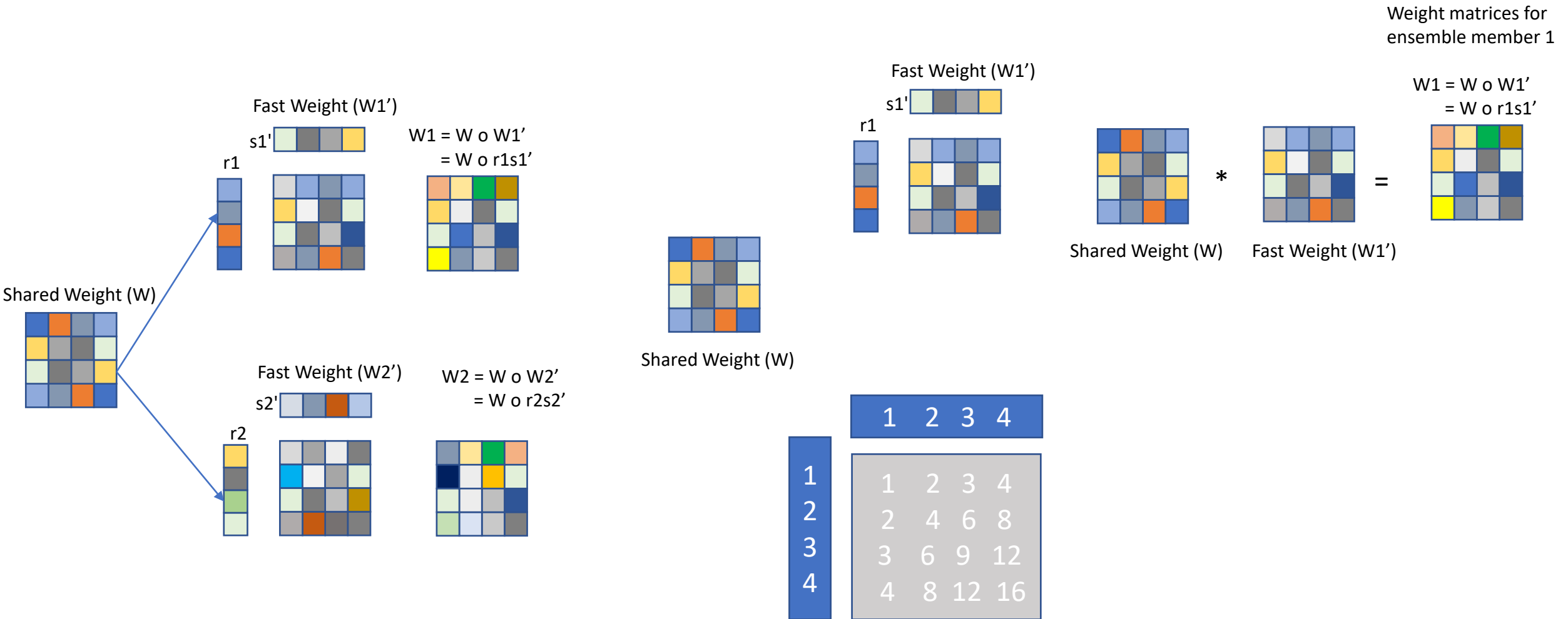
DNN (Single Model)



Deep Ensemble (size=2)

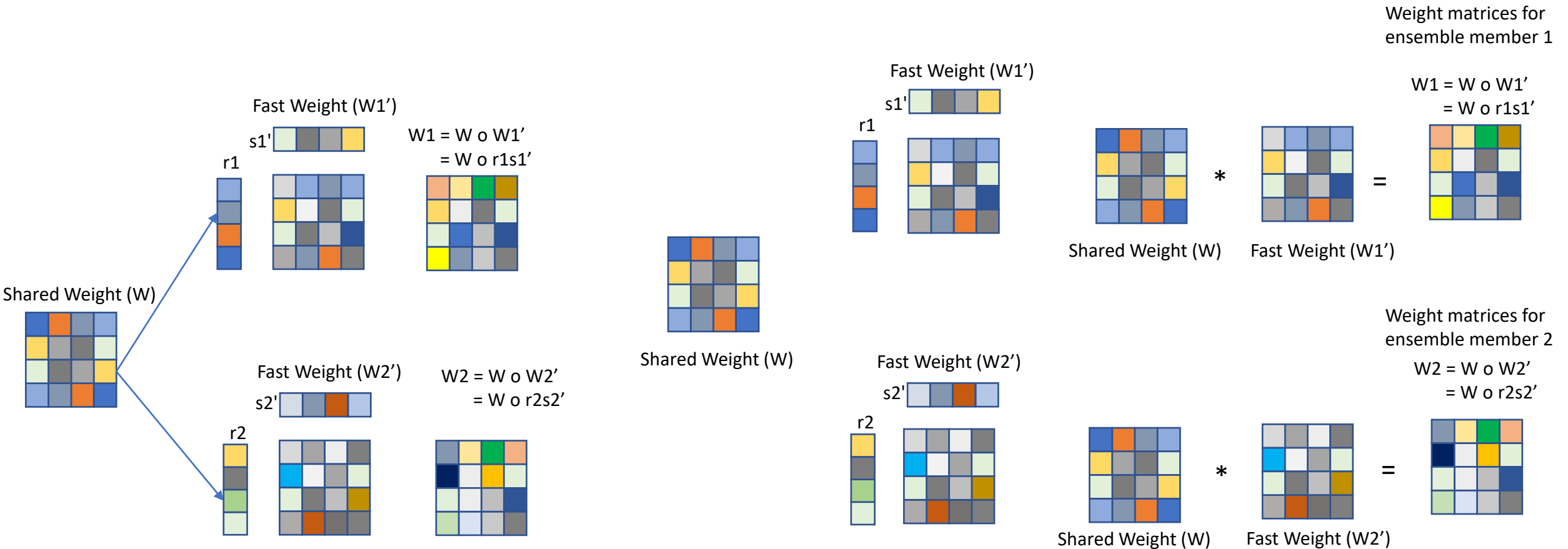


# Batch Ensemble (size=2)



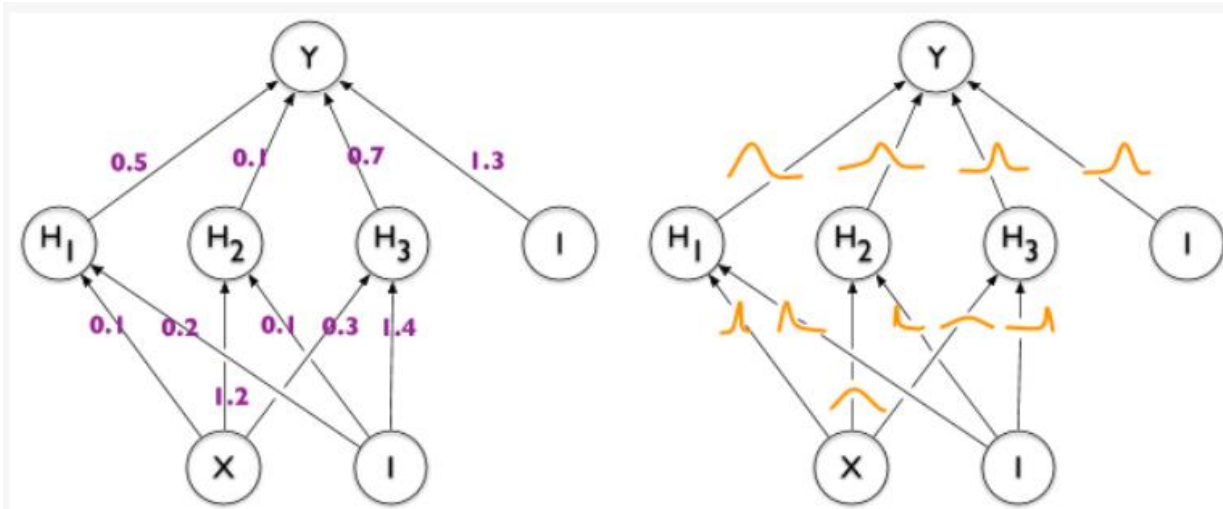
$r_i$  and  $s_i$  share same dimension as input and output –  $m$  &  $n$  respectively

# Batch Ensemble (size=2)

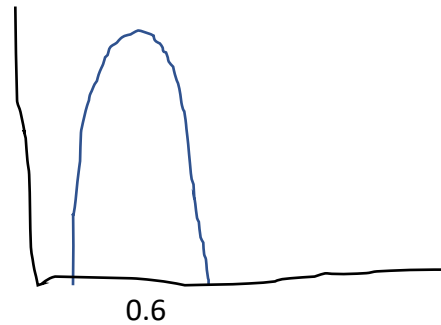
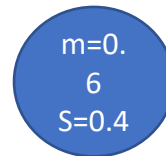


$r_i$  and  $s_i$  share same dimension as input and output –  $m$  &  $n$  respectively

# Rank-1 BNN (size=2)



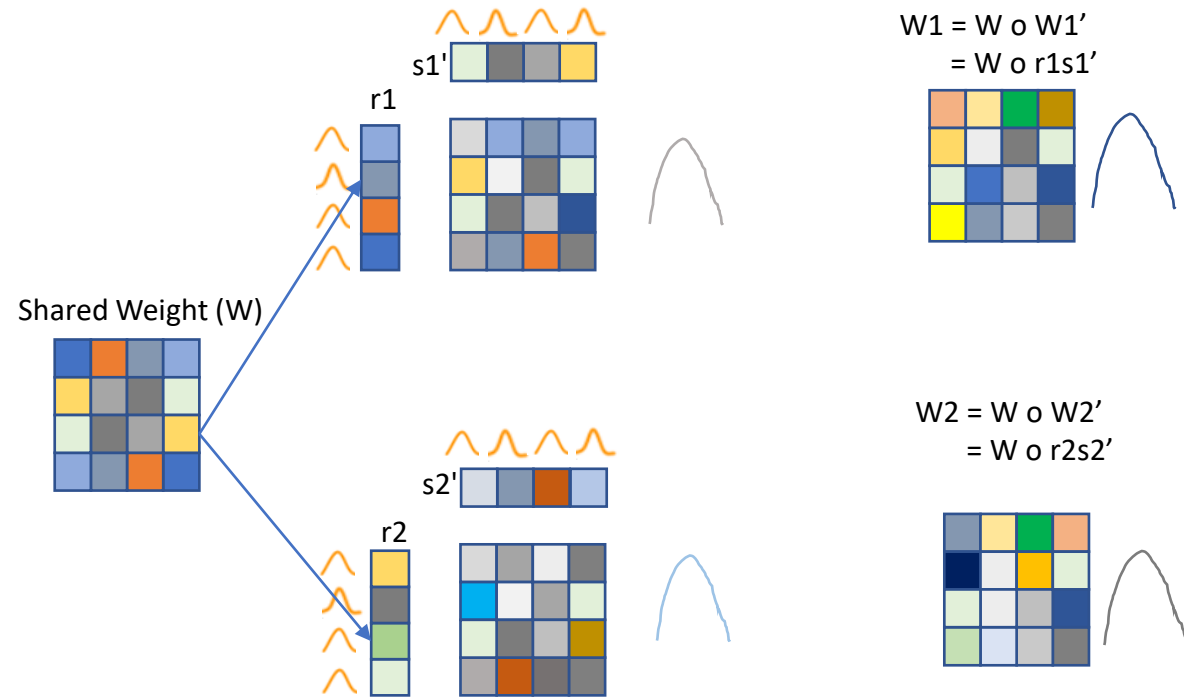
Unlike deterministic Neural Networks(left) that have a fixed value of their parameters, Bayesian Neural Networks(right) has a distribution defined over them.



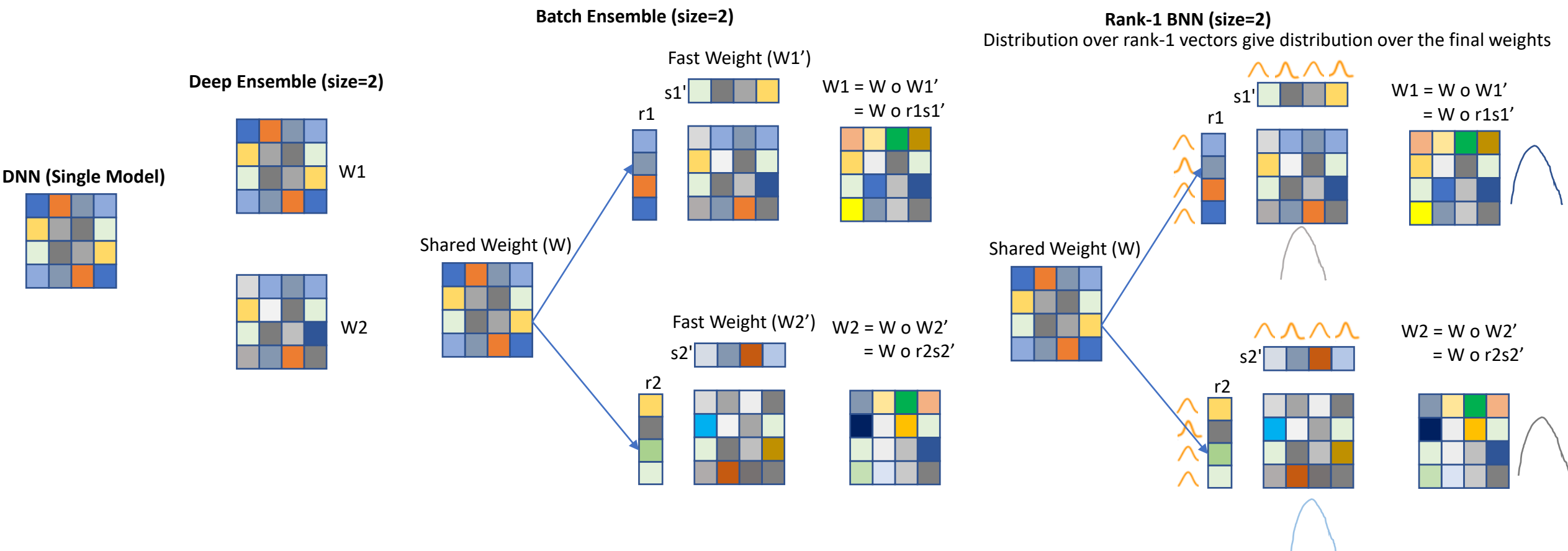


# Rank-1 BNN (size=2)

Distribution over rank-1 vectors give distribution over the final weights



# Deep Ensemble vs Batch Ensemble vs Rank-1 BNN

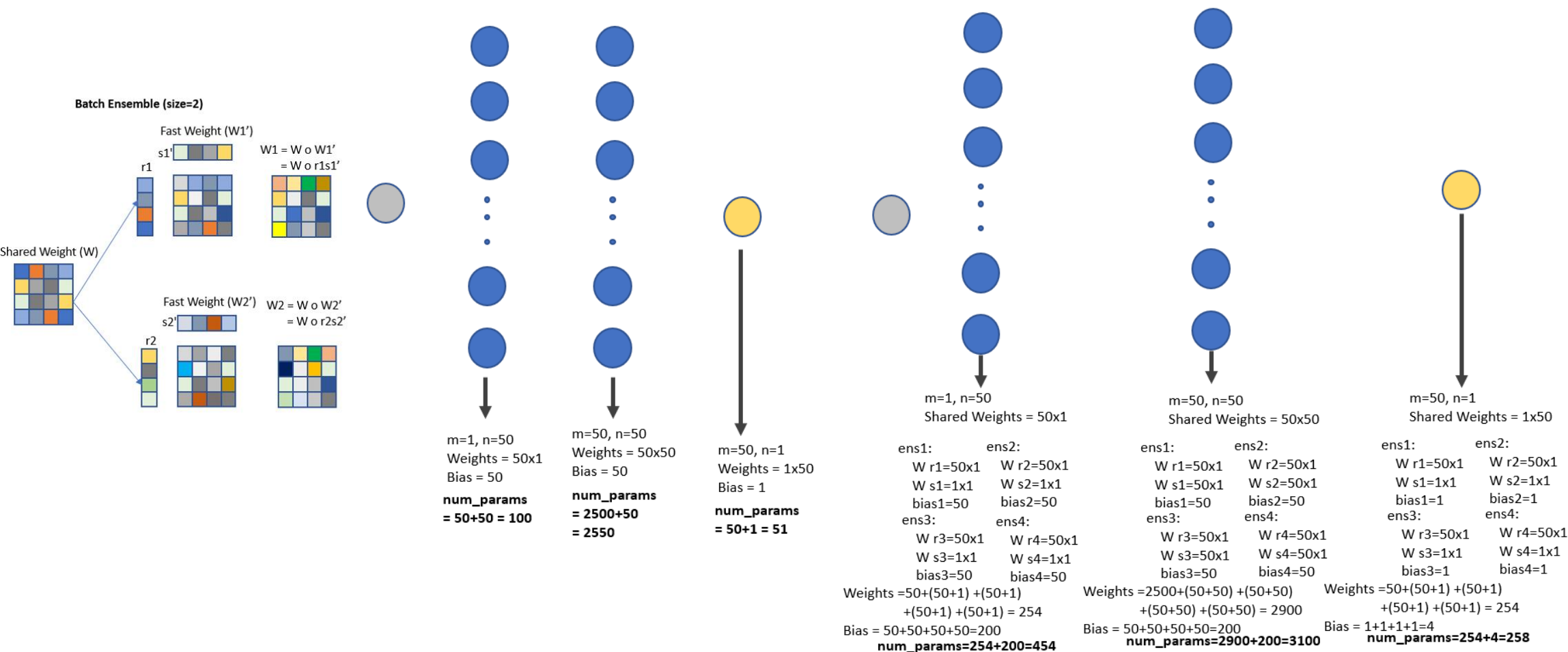


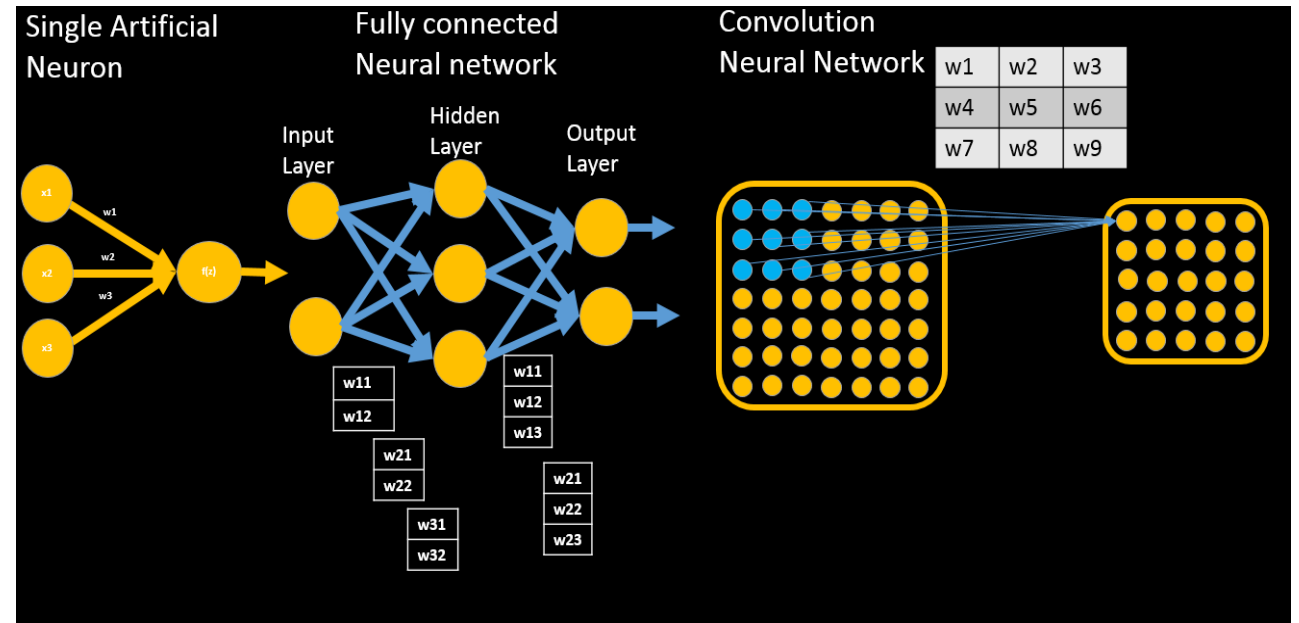
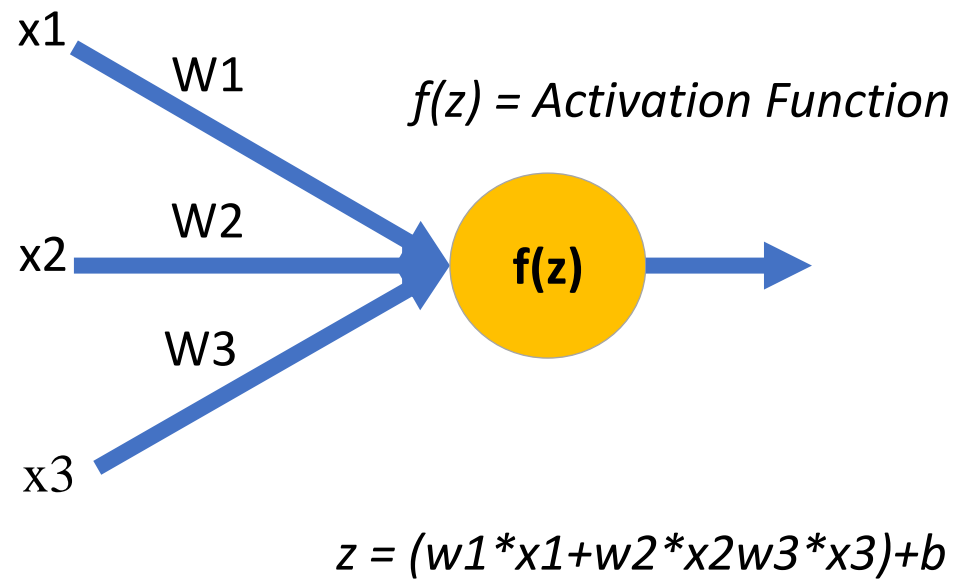
## Part-2

# Understanding **weight matrices** and **parameters count** of different deep learning ensemble models used for **Uncertainty Estimation**

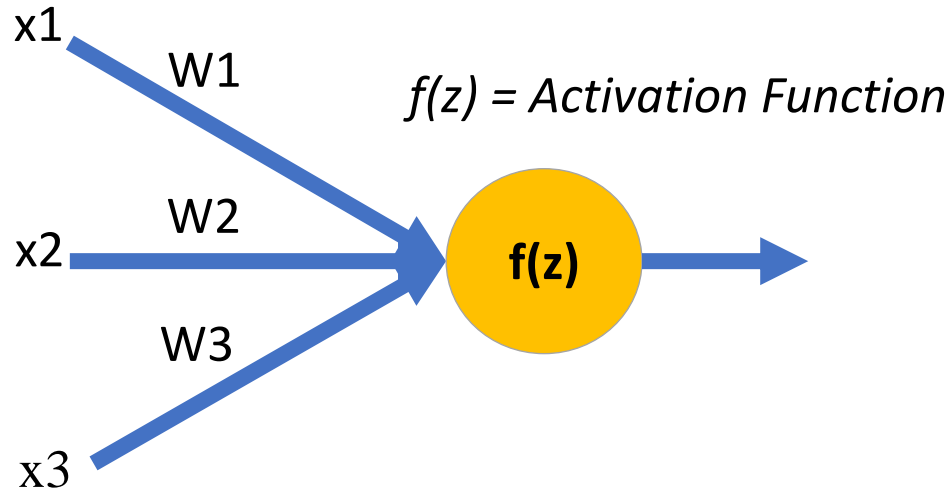
## Part-2: Code and Computation

### Deep Ensemble vs Batch Ensemble vs Rank-1 BNN





# Fully connected NN

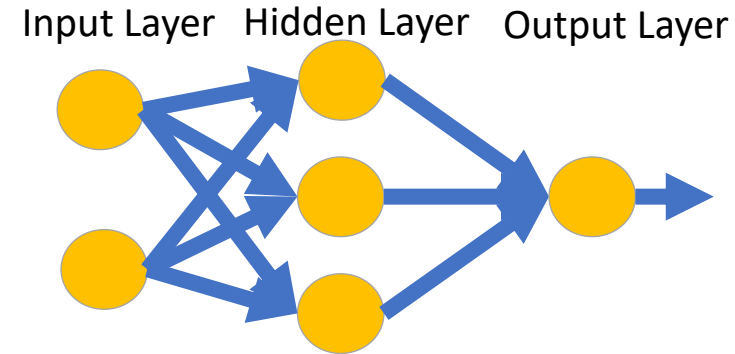


$$z = (w_1 * x_1 + w_2 * x_2 + w_3 * x_3) + b$$

Num\_params = weights + bias

For one output unit (neuron)

1. Weights = num of input unit
2. Bias = 1



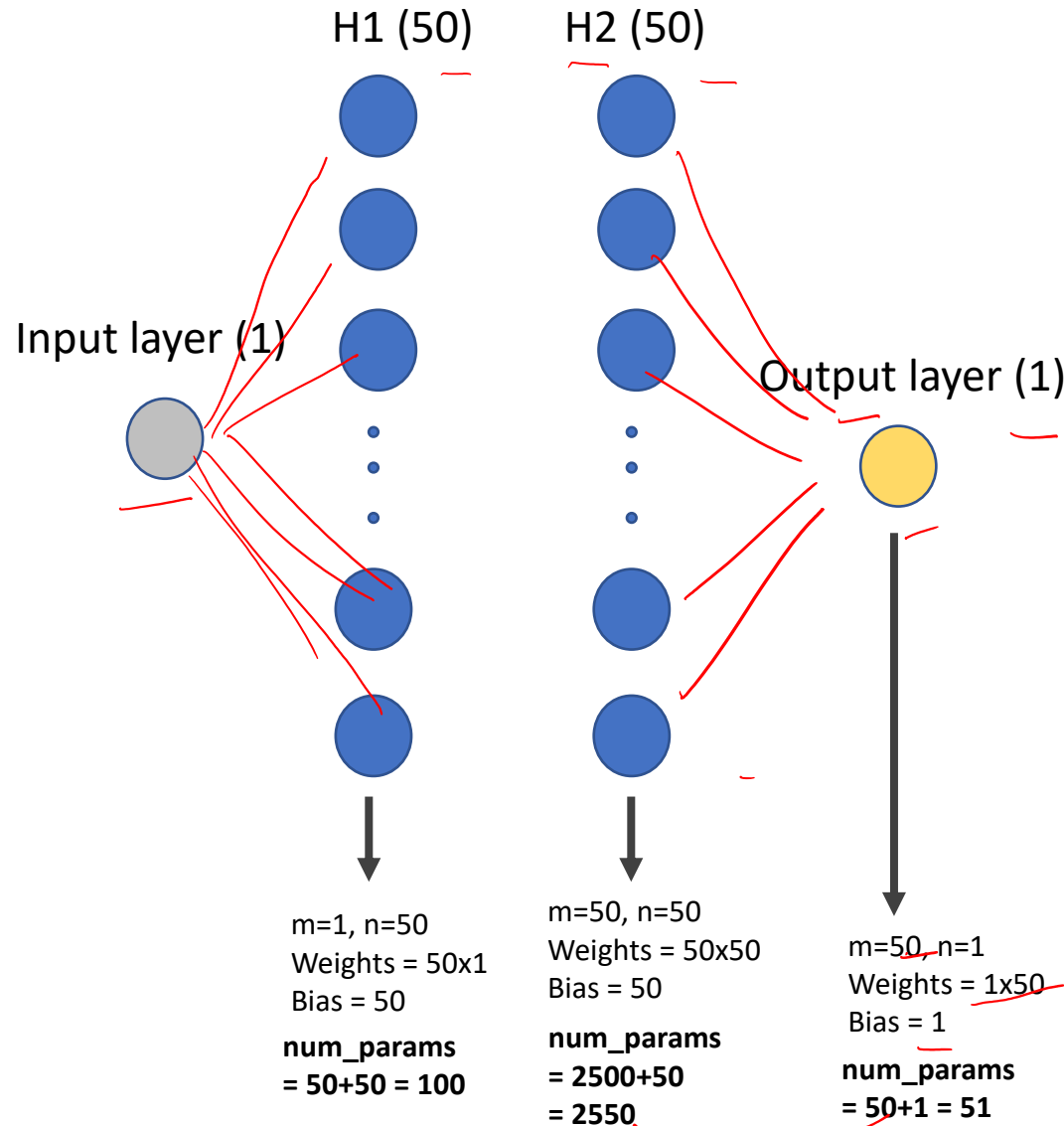
For any layer in NN

- $m$  = num input nodes
- $n$  = num output nodes

1. Weights =  $n * m$
2. Bias =  $n$

$$\text{Num\_param} = \text{weights} + \text{bias} = n * m + n$$

# Single model (Fully Connected NN)



For any layer in NN

- $m = \text{num input nodes}$
- $n = \text{num output nodes}$

1. Weights =  $n \times m$

2. Bias =  $n$

**Num\_param** = weights + bias =  $n \times m + n$

$$m=1$$

$$n=50$$

$$w = 50 \times 1$$

$$b = 50$$

$$\left. \begin{array}{l} w = 50 \times 1 \\ b = 50 \end{array} \right\} 50 + 50 = 100$$

$$m=50$$

$$n=50$$

$$w = 50 \times 50 = 2500$$

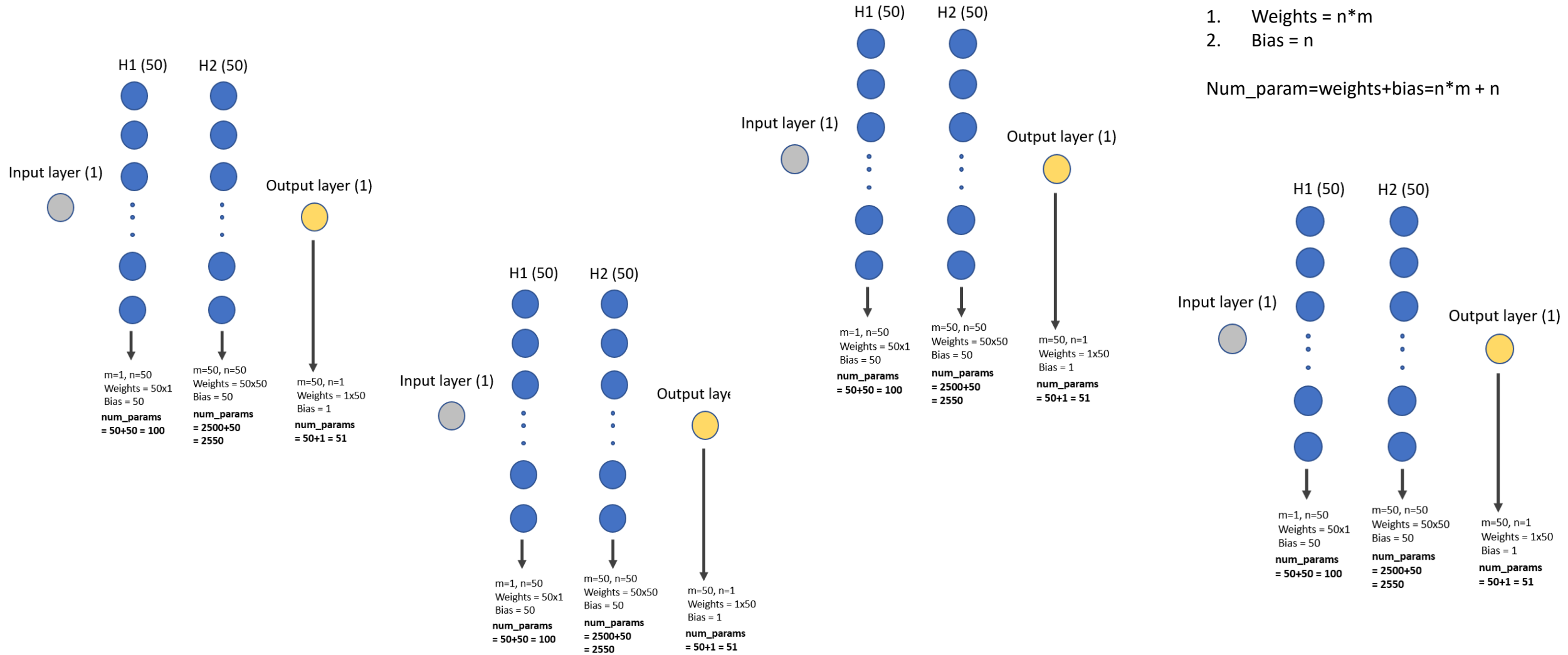
$$b = 50$$

$$\left. \begin{array}{l} w = 50 \times 50 = 2500 \\ b = 50 \end{array} \right\} 2550$$

$$2701$$

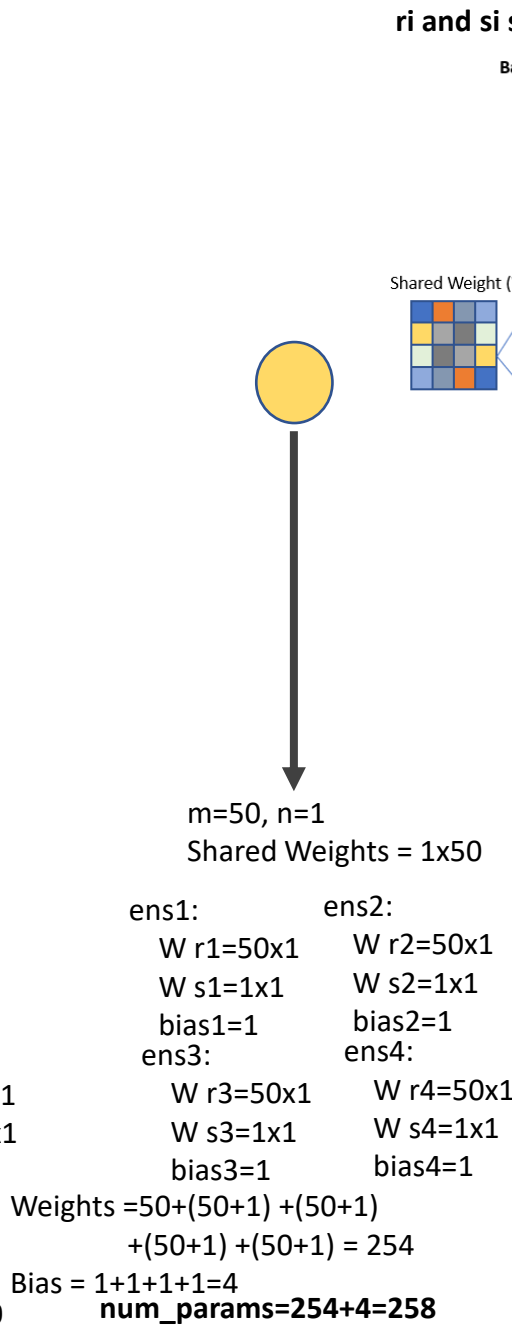
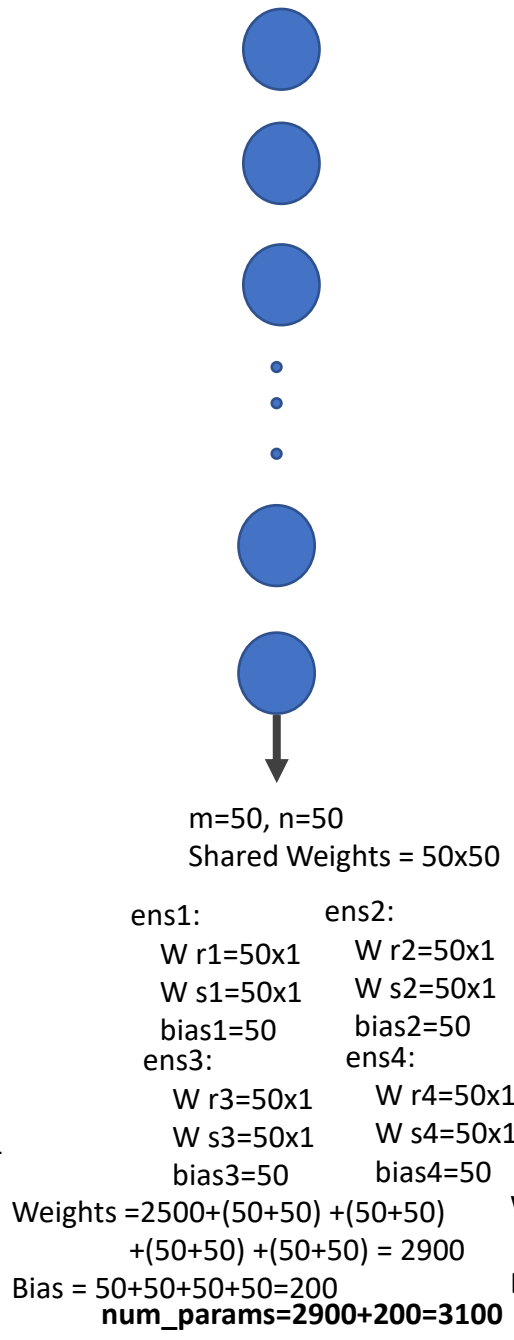
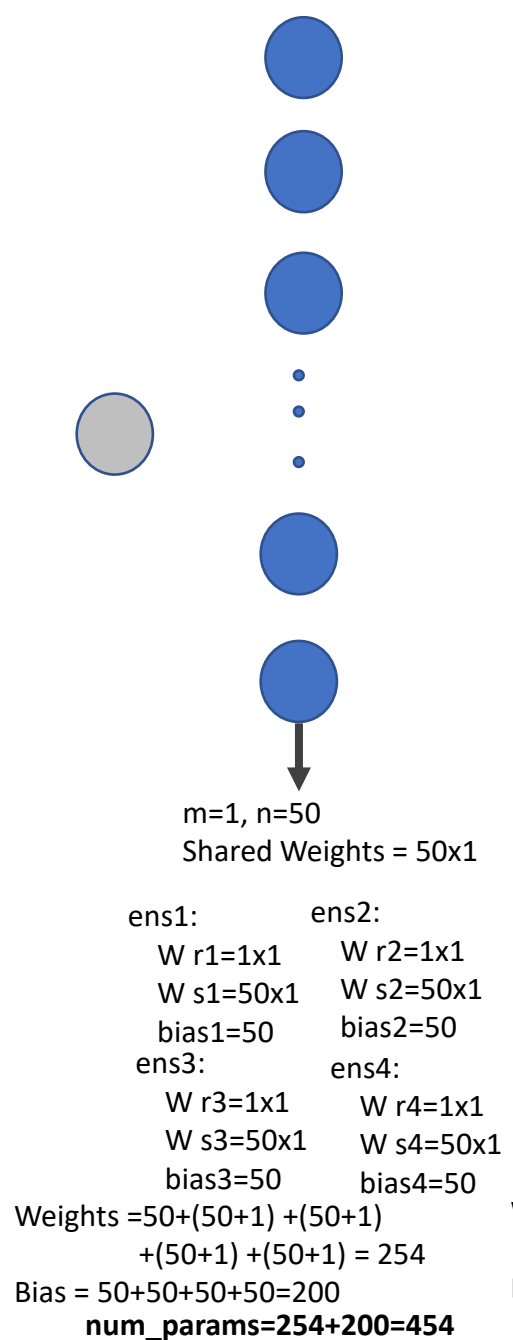
$$\begin{array}{r} 2550 \\ + 151 \\ \hline 2701 \end{array}$$

# Deep Ensemble (Ensemble size=4)

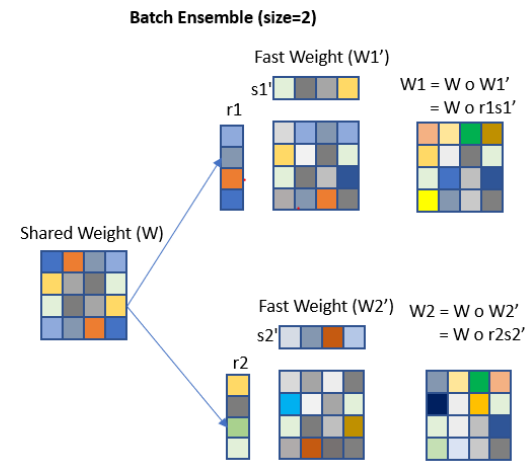




# Batch Ensemble (ensemble size = 4)



ri and si share same dimension as input and output – m & n respectively



For any layer in NN

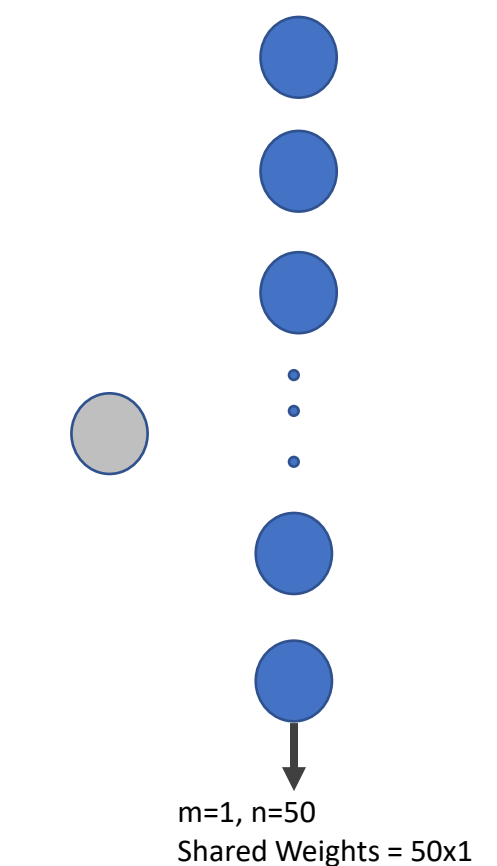
- $m = \text{num input nodes}$
- $n = \text{num output nodes}$

1. Weights =  $n * m$
2. Bias =  $n$

**Num\_param = weights + bias =  $n * m + n$**

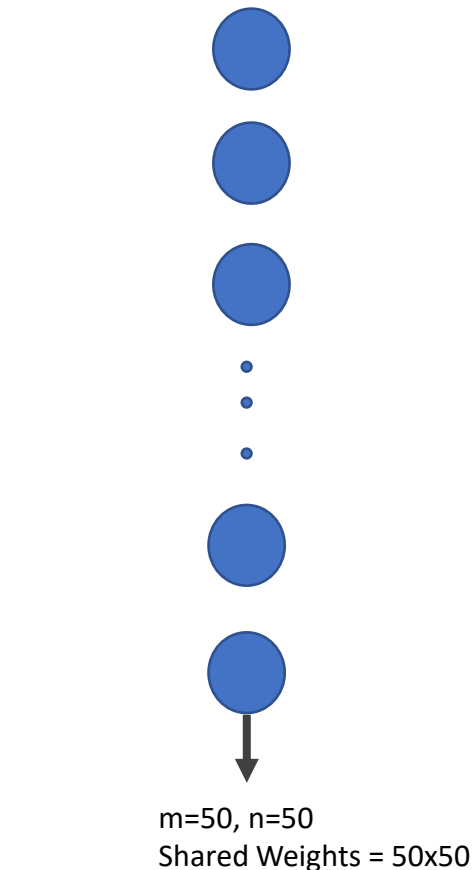


Rank-1 BNN for each of the rank1 vector we have a distribution that is represented by 2 value mu and sigma



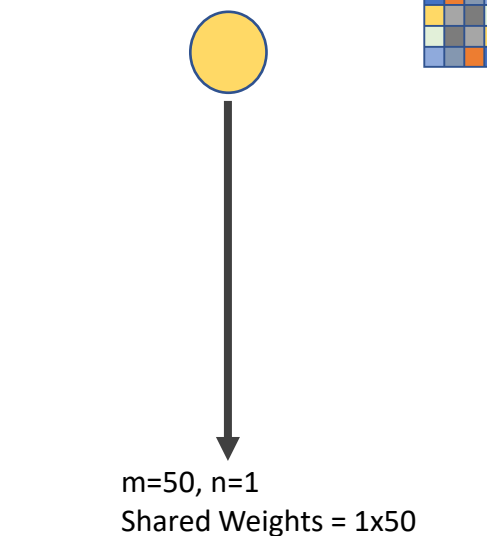
ens1:	ens2:
W r1=1x1	W r2=1x1
W s1=50x1	W s2=50x1
bias1=50	bias2=50
ens3:	ens4:
W r3=1x1	W r4=1x1
W s3=50x1	W s4=50x1
bias3=50	bias4=50

Weights =  $50 + (50+1)*2 + (50+1)*2$   
 $+ (50+1)*2 + (50+1)*2 = 458$   
 Bias =  $50+50+50+50=200$   
**num\_params=458+200=658**



ens1:	ens2:
W r1=50x1	W r2=50x1
W s1=50x1	W s2=50x1
bias1=50	bias2=50
ens3:	ens4:
W r3=50x1	W r4=50x1
W s3=50x1	W s4=50x1
bias3=50	bias4=50

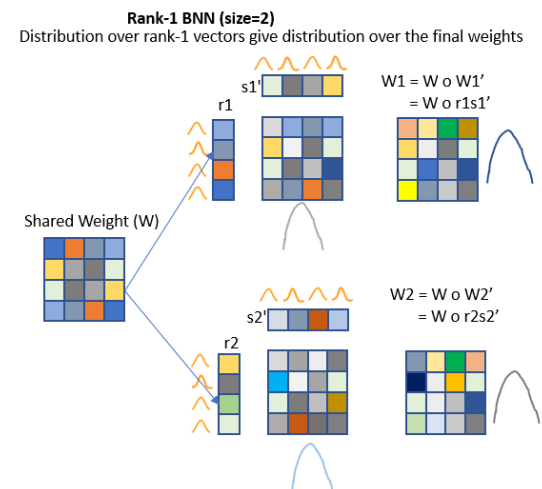
Weights =  $2500 + (50+50)*2 + (50+50)*2$   
 $+ (50+50)*2 + (50+50)*2 = 3300$   
 Bias =  $50+50+50+50=200$   
**num\_params=3300+200=3500**



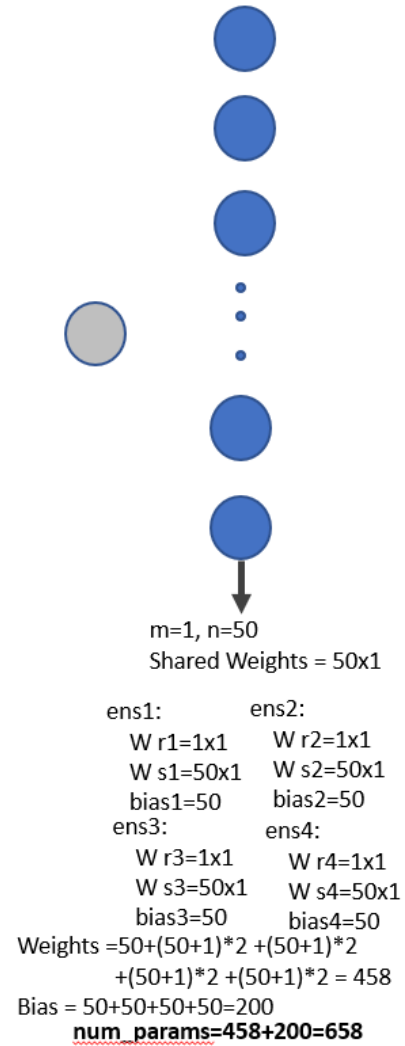
ens1:	ens2:
W r1=50x1	W r2=50x1
W s1=1x1	W s2=1x1
bias1=1	bias2=1
ens3:	ens4:
W r3=50x1	W r4=50x1
W s3=1x1	W s4=1x1
bias3=1	bias4=1

Weights =  $50 + (50+1)*2 + (50+1)*2$   
 $+ (50+1)*2 + (50+1)*2 = 458$   
 Bias =  $1+1+1+1=4$   
**num\_params=458+4=462**

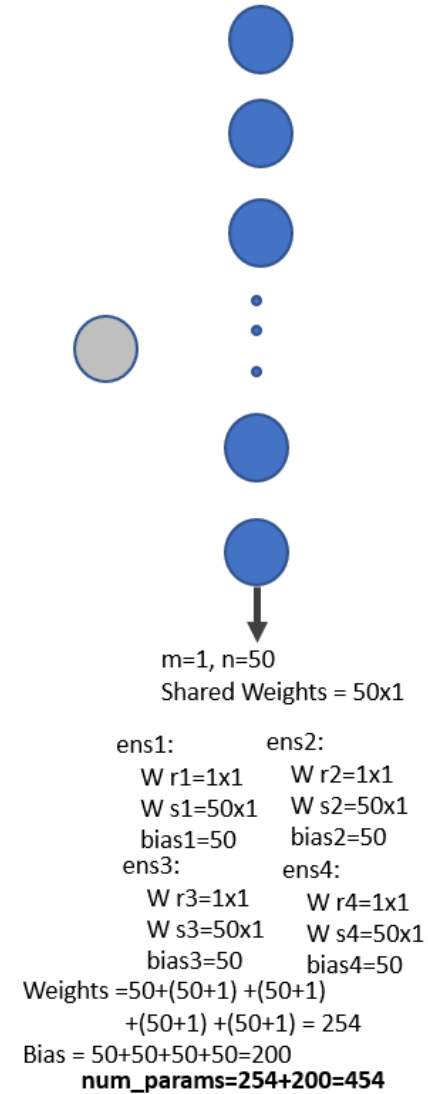
ri and si share same dimension as input and output – m & n respectively



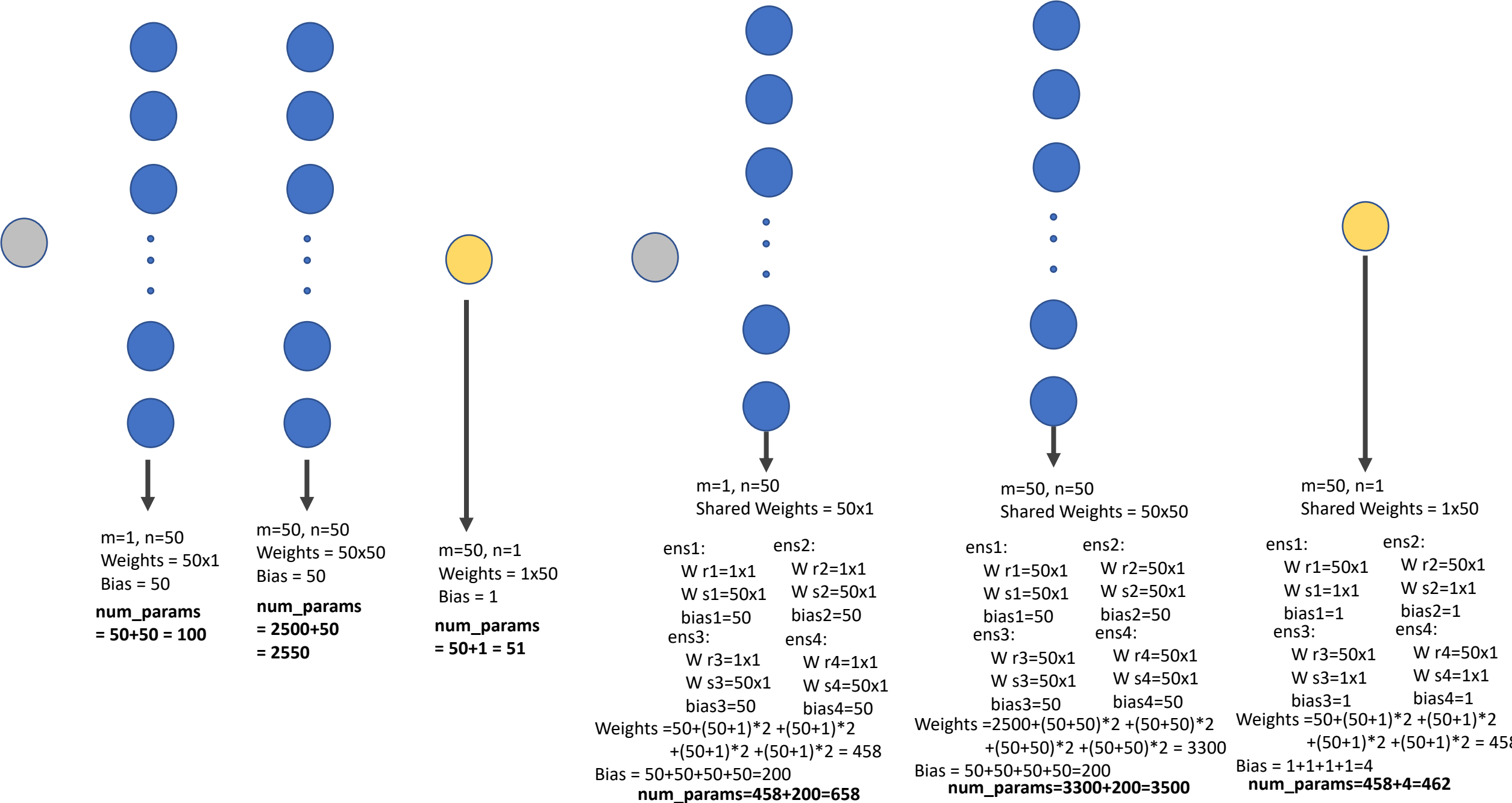
## Rank-1



## Batch Ensemble



Rank-1 BNN for each of the rank1 vector we have a distribution that is represented by 2 value mu and sigma



Thank You