

## II Element of Hypothesis Test.

### B. Bayesian Hypotheses Testing.

$$H_0: Y \sim P_0 \quad \text{vs} \quad H_1: Y \sim P_1$$

$Y$  has distribution  $P$ .

$H_0$ : Null,  $H_1$ : alternative.

Decision rule  $\delta$  on  $\Gamma$

$$\delta(y) = \begin{cases} 1 & \text{if } y \in \Gamma_1 \\ 0 & \text{if } y \in \Gamma_1^c. \end{cases}$$

$$C_{2y}, \{y=0, 1\} \text{ and } g=0, 1$$

$C_{1y}$ : the cost of choosing

$H_1$  when  $H_1$  is true.

$$R_g(\delta) = C_{1y}P_g(\Gamma_1) + C_{0y}P_g(\Gamma_0), g=0, 1$$

$$R_g(\delta) =$$

cost of choosing  $H_1$  when  $H_1$  is true

$\times$   
The probability of to do so.

cost of choosing  $H_0$  when  $H_1$  is true.

$\times$   
The probability of to do so.

$\pi_0$ : the prior probability of occurrences of  $H_0$

$\pi_1$ :  $\sim$  of  $H_1$

Average of Bayes risk:

$$r(\delta) = \pi_0 R_0(\delta) + \pi_1 R_1(\delta)$$

$$R(\delta) = \sum_{j=0}^1 \pi_j C_{0j} + \int_{P_1} \left[ \sum_{j=0}^1 \pi_j (C_{1j} - C_{0j}) p_s(y) \right] \mu(dy)$$

thus. minimum over  $P_1$  if we choose.

$$\Gamma_1 = \{y \in P \mid \pi_1(C_{11} - C_{01}) p_1(y) \leq \pi_0(C_{00} - C_{10}) p_0(y)\}$$

Assume  $C_{11} < C_{01}$

$$\Gamma_1 = \{y \in P \mid p_1(y) \geq \tau p_0(y)\}, \quad \tau = \frac{\pi_0(C_{10} - C_{00})}{\pi_1(C_{01} - C_{11})}$$

Likelihood-ratio test.

$$L(y) = \frac{p_1(y)}{p_0(y)} \quad (\text{likelihood ratio})$$

$$S_B(y) = \begin{cases} 1 & \text{if } L(y) \geq \tau \\ 0 & \text{if } L(y) < \tau \end{cases}$$

$\pi_0(y), \pi_1(y)$  are posterior of two hypotheses.

$$\Gamma_1 = \{y \in P \mid (C_{10}\pi_0(y) + C_{01}\pi_1(y)) \leq (C_{00}\pi_0(y) + C_{11}\pi_1(y))\}$$

$(C_{10}\pi_0(y) + C_{01}\pi_1(y))$  posterior cost of choosing  $H_1$  when  $y$  is observed.



Example:

$Y$  (output of channel) can be 1, 0.

True msg	received	chance
Zero	1	$\lambda_0$
Zero	0	$1 - \lambda_0$
One	0	$\lambda_1$
One	1	$1 - \lambda_1$

$\lambda_j$ ,  $j$  was transmitted

$T$  is  $\{0, 1\}$  observation set.

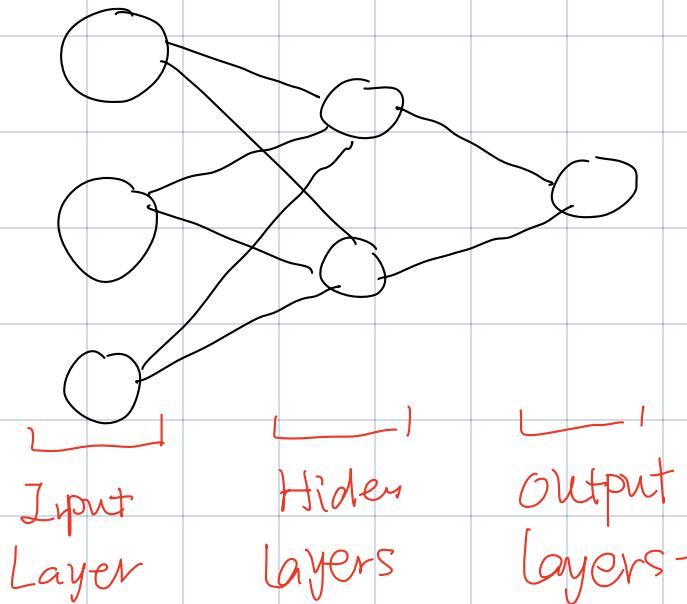
$Y$  has densities  $p_j(y) = \begin{cases} \lambda_j & \text{if } y \neq j \\ 1 - \lambda_j & \text{if } y = j \end{cases}$

24, Oct, 2024 WP 1.2

## Deep Learning Crash Course for Beginners (YouTube Tutorial)

DL can learn features from raw data.

A simple Neural Network



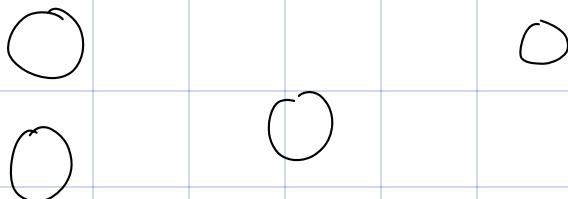
Learning process

Forward Propagation:



$$\hat{y} = \sigma \left( \sum_{i=1}^n x_i w_i + b_j \right)$$

$\sigma$ : activation function -



## Back Propagation.

$\hat{y}$  - Right

Wrong  $\rightarrow$  Loss function quantify the deviation from expected output.

$\rightarrow$  Go backwards & adjust initial weights and bias.

## Terminology -

### Activation Functions:

- Introduce Non-Linearity.
- Decide whether a neuron can contribute to next layer.

Sigmoid Function. ( $\text{sig}(t) = \frac{1}{1+e^{-t}}$ )

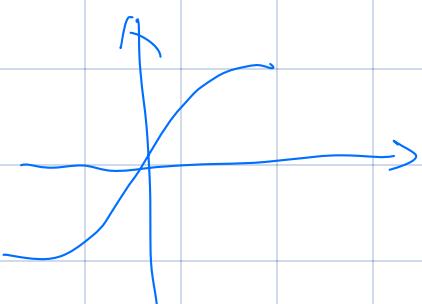


Non-linear, Analog outputs.

$$\text{sig}(t) \in (0, 1)$$

gradient vanish - problem.

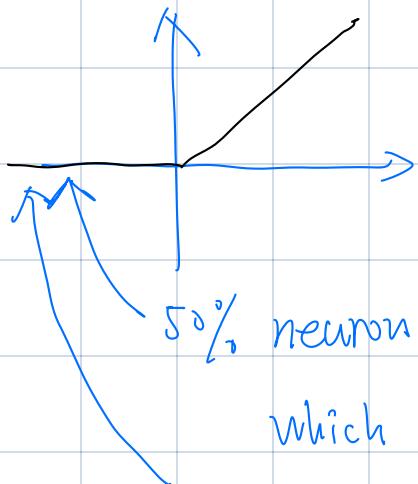
Tanh Function.  $(\frac{e^x - e^{-x}}{e^x + e^{-x}} - 1) = \tanh(x) = 2\text{sigmoid}(2x) - 1$



Derivative steeper than Sigmoid.

still gradient vanish problem -

# ReLU Function. (Rectified Linear Unit)



sparse Activations.

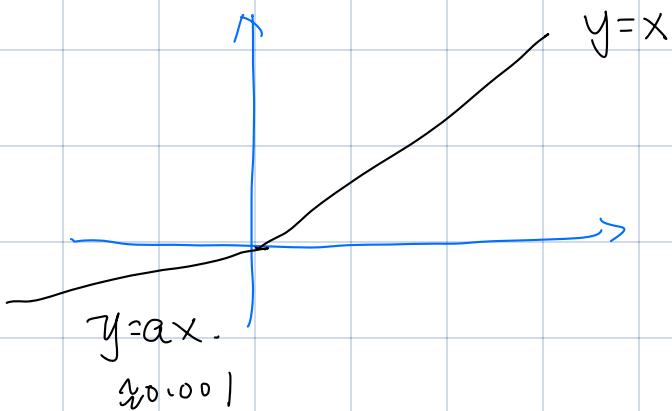
50% neuron initialized at negative axis.

which make network lighter.

Gradient = 0, Dying ReLU Problem.

Not get adjusted in descent.

Leaky ReLU Function.



ReLU is less computationally expensive than Tanh or sigmoid.

Loss Function.

Regression, Binary Classification, Multi-Class Classification.

Optimizers.

minimize loss function.

by updating the network based on the output of Loss function.

Loss function guide optimizers.

## Gradient Descent.

$$\nabla f(x, y) = \left( \frac{\partial f}{\partial x}(x, y), \frac{\partial f}{\partial y}(x, y) \right)$$

Gradient always points in the steepest direction.

Choose proper Learning Rate to avoid local minimum.

too large — overshoot

too small — local minimum.

## Stochastic Gradient Descent.

Using subset of sample "batch"

- less computation.

using momentum.

- get rid of local minimum

$$\theta_{t+1} = \theta_t - \eta \nabla J(\theta_t)$$

$$v_{t+1} = \beta v_t + \nabla J(\theta_t)$$

$$\Rightarrow \theta_{t+1} = \theta_t - \eta v_{t+1}$$

## Adagrad:

adapts learning Rate to individual features.

Ideal to sparse dataset.

but:

$\eta$  get small overtime.

$$c_{t,i} = c_{t-1,i} + (\nabla J(\theta_t)_i)^2$$

accumulating gradient<sup>2</sup>

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{c_{t,i} + \epsilon}} \nabla J(\theta_t)_i$$

## RMSprop

accumulates Gradients in a fixed window.

$$E[g^2]_t = \beta E[g^2]_{t-1} + (1-\beta)(\nabla J(\theta_t))^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} \nabla J(\theta_t)$$

27 Oct. 2024 WPI.2

## Adam ( Adaptive Moment Estimation )

Using momentum too

Pretty widespread today.

## 4. Parameters & Hyperparameters.

### Model Parameter.

Value define the skill of the model

estimated directly from data.

like: weight, bias.

### Model Hyperparam.

set manually to control the behavior of model

like: learning rate, layer and neuron number

## Epochs, Batches, Batch size and iterations.

epochs: entire dataset pass to network;

too many epochs lead to overfitting.

batches: subset of training dataset used to model in one iteration.

batch size: total number of training examples in a batch.

iterations: number of batches need to complete one epoch.

Introduce to learning:

Supervised Learning:

Learn by exampl.

train model on well-labelled data.

- Predict the correct label for unseen data.

sub category.

Classification?

creating mapping function.

i Handwriting Digits Dataset

Algorithms?

Linear classifier, Support vector machines,

K-Nearest Neighbour, Random forest.

Regression:

find a relationship between dependent & independent variables.

Algorithms:

Linear regression, Lasso regression, Multivariate regression.

## Unsupervised learning:

find features of data; no label.  
hidden pattern

subcategory?

clustering:

grouping data into different clusters of groups

Partitional clustering:

a data point can belong to a single cluster

hierarchical clustering:

clusters Within clusters

data point may belong to many clusters.

Algorithms:

K-Means, expectation maximization, HCA

28. Oct. 2024. WPl.2.

Association:

find relationships between different entities.

Reinforcement learning: (goal-oriented)

enable an agent to learn in interactive environment  
by trial & error,

Rewards & punishment used to guide agent

Maximize total cumulative reward.

## Regularization.

Core problem: overfitting.

performing well on training set  
but not on testing set.

### Tackling Overfitting.

Dropout: randomly removes some nodes.

Augmentation: fake data; More Data, Better model.

early stopping: stop when validation error increases.

## Neural Network Arch.

### Fully-Connected Feed Forward Neural Networks (Vanilla)

each neuron is connected to every subsequent layer with no backward connections.

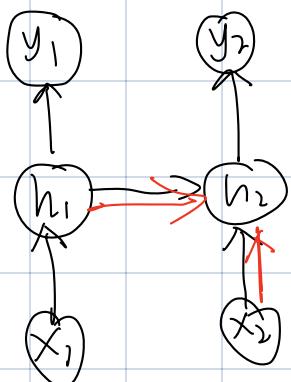
[Inputs, outputs, Hidden Layers, Neurons per hidden layer, activation functions]

⇒ wider array of problems.

- Can't handle sequential data. (sentence...)
- don't share parameters across time.

# Recurrent Neural Net. (RNN)

use feedback loop in hidden layers.



use previous states as input  
for its current prediction.

use backpropagation algorithm.

## Vanishing Gradient Problem.

gradient will shrink over network due to  
the nature of backpropagation, leading slow learning  
of front layers.  $\Rightarrow$  short-term memory-

## Solutions:

Long short term memory (LSTM)

Gated RNN (GRNN)

CRNN has two gates:

Update gates -

Reset gates -

LSTM has three gates:

update gates, Reset Gate, Forget gate

# Convolutional NNs -

good for processing data like:  
Images, audio and video

## Hidden Layers in a CNN

Convolution layer -

convolution and pooling  
function are used -

Pooling ~

Fully-Connected ~

Normalization ~

## Convolution

extract visual features in chunks (filter  
or kernel)

## Pooling -

reduce the number of neurons necessary in  
subsequent layers.

Max pooling : pick biggest number in selected  
region.

Min pooling : ~

# Creating a deep learning model

## Gathering Data.

right data is key.

size:

Amount of data  $\approx 10 \times \frac{\text{No. of model}}{\text{parameters}}$

Quality:

archive.ics.uci.edu-

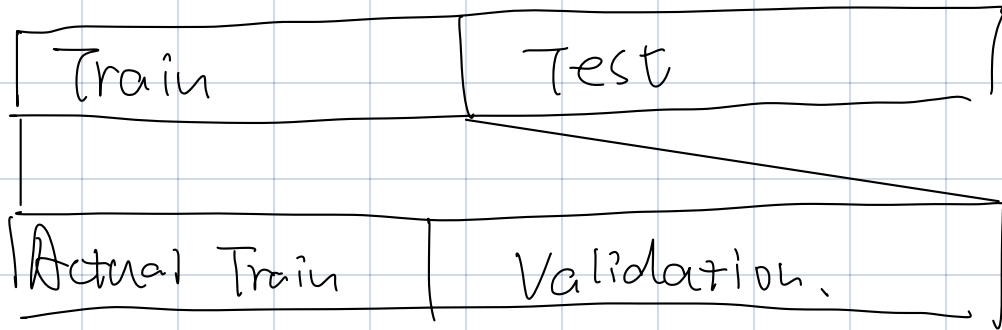
Kaggle.com/datasets-

Google Dataset Search.

## Preprocessing the Data.

o splitting dataset into subsets.

: training, validation, testing.



Cross-Validation.

o formatting:

might not in the format we need.

## o Missing Data.

represented as "NaN" or "Null"

1. eliminating features with missing values.

but risk losing relevant info.

2. inputting the missing values.

## o Sampling.

use smaller sample if large dataset,

imbalancing data:

learning majority class > learning minority class

$\Rightarrow$  bias to majority.

Example weight = original weight  $\times$  DownSampling Factor -

$\Rightarrow$  increase weight of majority to compensate decrease in number of samples.

## o Feature Scaling.

1. Normalization.

2. Standardization. ~~X~~

Train.

Evaluation -

test model on the validation set.

## Optimizing -

tuning Hyperparameter-

No. epochs  $\uparrow$

adjust learning rate -

addressing Overfitting -

- o get more data.

- o Reduce model size.

- o Regularization.

$L_1$  Reg.  $\Rightarrow$  add cost to limit complexity

$L_2$  Reg of model.

- o Data augmentation -

In case of image - [flipping, blurring, zooming]

- o Dropout.

## 14. Nov. 2024.

--model, --channel, --src, --dst, --shr, --member of set

① training set [0, 1] or some other value

② AWGN - SNR

Rician fading -  $K$ : ratio between power in direct path.

and that in other

$\rightarrow$   $\Omega$ : total power from both paths.

Rayleigh fading -

Generate noised signal [selecting channel type]

Selecting Model, channel type, sizeof DataSet -

generating training set: output a files

train

$0.01 \sim 0.10$

Adam

channel-type

Model-type

BER

AWGN

DMN

0.0883

1.00~

0.0387

CNN

0.1010

0.0733

0.1258

0.1104

SGD

# of Layers 0.0248

Optimize, loss function, activation 0.02150

Epoch = 20? Does matrix matter?

0.02950

What channel I record for log

0.002 Rayleigh  
0.00197 Rician

## Meeting 20 Nov.

- Using sigmoid at output layer.
- Noised signal should be complex number.
- Try to implement first-order mate-learn.

Outcome after solve complex and sigmoid.

Epochs = 50, batch size = 32, data size = 10k  
test size = 0.2 × 10k.

512, leaky\_relu.

256, ~

128, ~

2, sigmoid.

Channel BER

AWGN 0.03187

Rayleigh 0.14497

Rician. 0.13916

Jan. 30, 2025.

Cognition refreshing.

train DL1 with AWGN normally.

Test DL1 with AWGN test set.

train DL2 with Rayleigh - ..

- - -

train DL3 with Rician.

- - -

train model4 with 3 noise types.

11/2/2025-

Failed to learn Rayleigh.

Overfit the AWGN

D1	AWGN ,	} Conventional DNN
D2	Rician ,	
D3	Rayleigh ,	
D4	Random Mixed channel ,	
D5	Sequential Mixed channel,	

D6	Random Mixed channel ,	} Meta-learning DNN
D7	Sequential Mixed channel,	

## Generalization Testing

D<sub>1~7</sub> 3GPP with limited data.

Control Both # of Train Data.

and # of Validating Data.

Train Date analogizes data collected from environment. for training

Validating Date analogizes data inputed after

2025.Feb.21

Project flow built:

Train Phase:

D1 AWGN,

D2 Rician,

D3 Rayleigh,

D4 Random Mixed channel,

D5 Sequential Mixed channel,

} Conventional  
DNN

D6 Task  
| Random\_mixed  
| Sequential\_mixed  
→ Meta-learning  
DNN  
with Retile  
Algorithm.

Fine Tune Phase:

D1-6 3CPP [100, 500, 1000]

Generalizing Testing Phase:

Validated after each Fine-tune

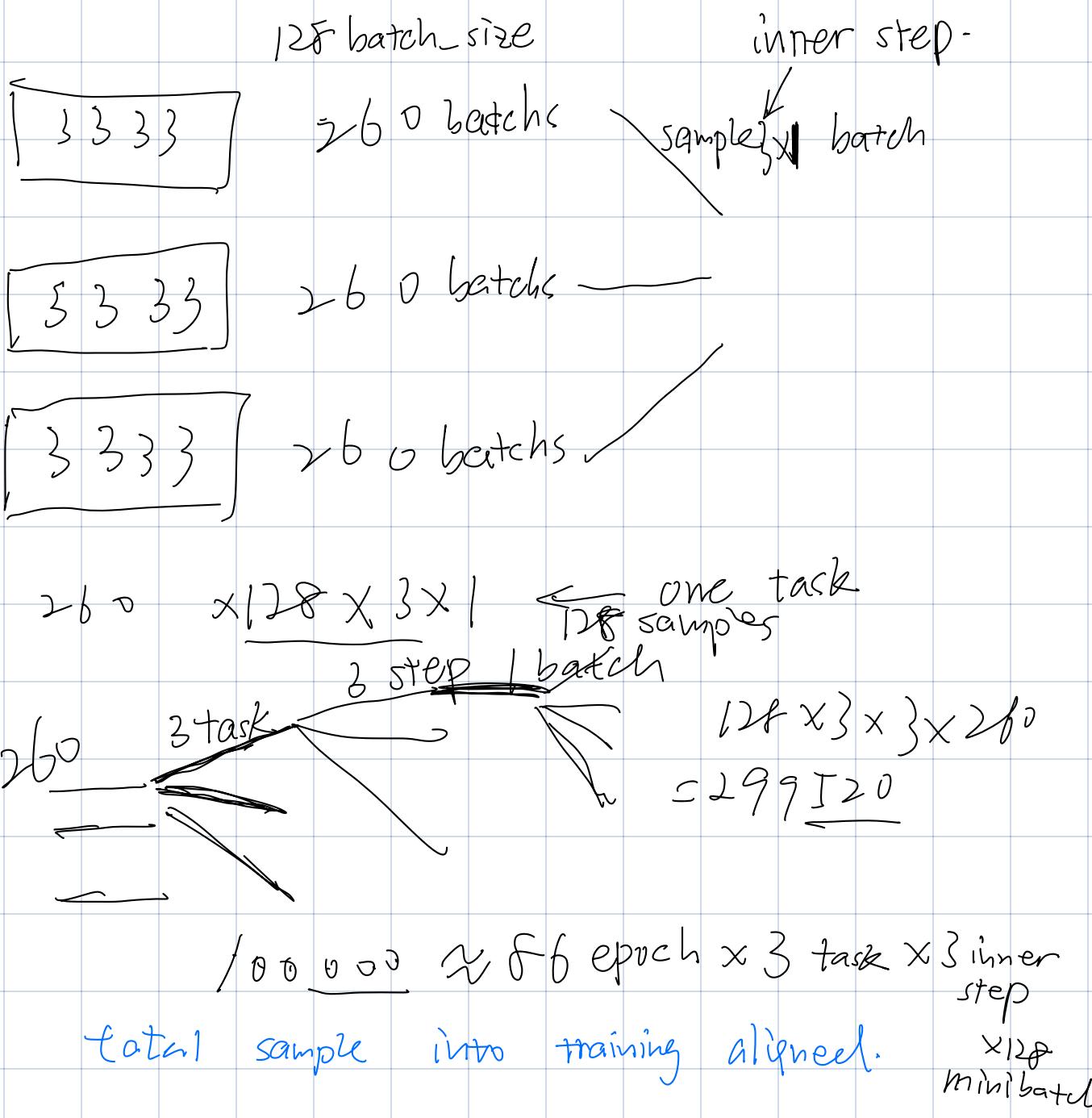
Loss and BER of D1 (AWGN) get lowest point at  
2nd, 3rd Epoch but bounce back in following Epoch.

→ learn rate too large?

D<sub>2-5</sub> present decreasing trend in both loss and BER.

→ learn rate too small / more epoch needed.

In generalization test, meta learning appears learning nothing (BER ≈ 0.5)



2025/2/27

## Reptile from OpenAI.

introduce meta-batch

与其每 epoch 都会在 3 个 task 学习

meta-batch 将 model 在不同数量 task 上

学习

Meta-learning descent too slow.

meta epoch : 10  $\Rightarrow$  2000

inner learn rate : 0.001  $\Rightarrow$  0.003

reptile meta learning rate : 0.01  $\Rightarrow$  0.25

↓

$$\text{self-adjust new\_meta-lr} = \left(1 - \frac{\text{epoch}}{\text{total epoch}}\right) \times \text{meta-lr.}$$

The BER of meta DNN decreased on the scale of 2000 iteration.

But still worse than Traditional DNN

Meta iteration  $\neq$  DNN epoch.

$\Rightarrow$  Number of times parameter updated.

For 100k sample to Traditional DNN

$$\# \text{Update} = \frac{\text{Total Sample}}{\text{batch\_size}} \times \text{Number of Epoch.}$$

Reptile

# update = meta-iteration.

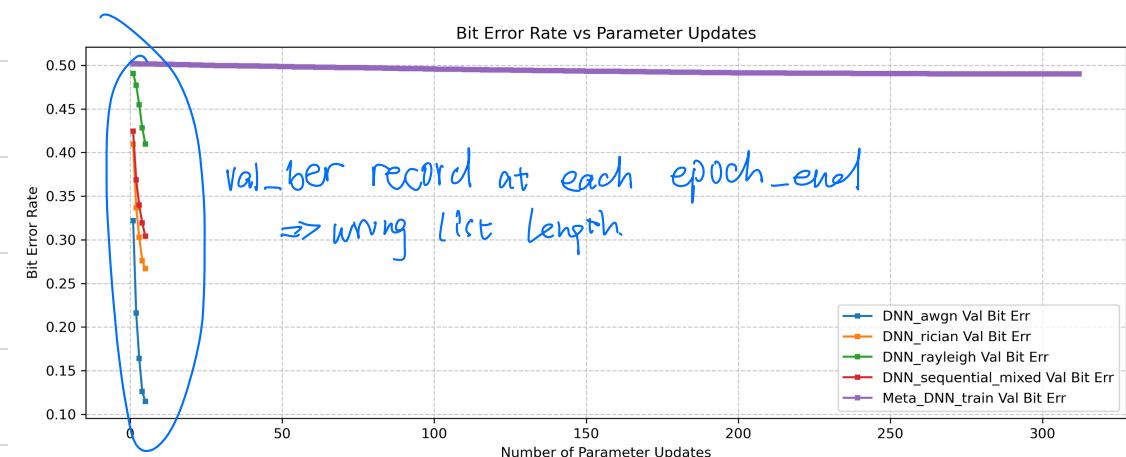
passed a list into a list  $\Rightarrow$  failed to output figure.

2025/3/4

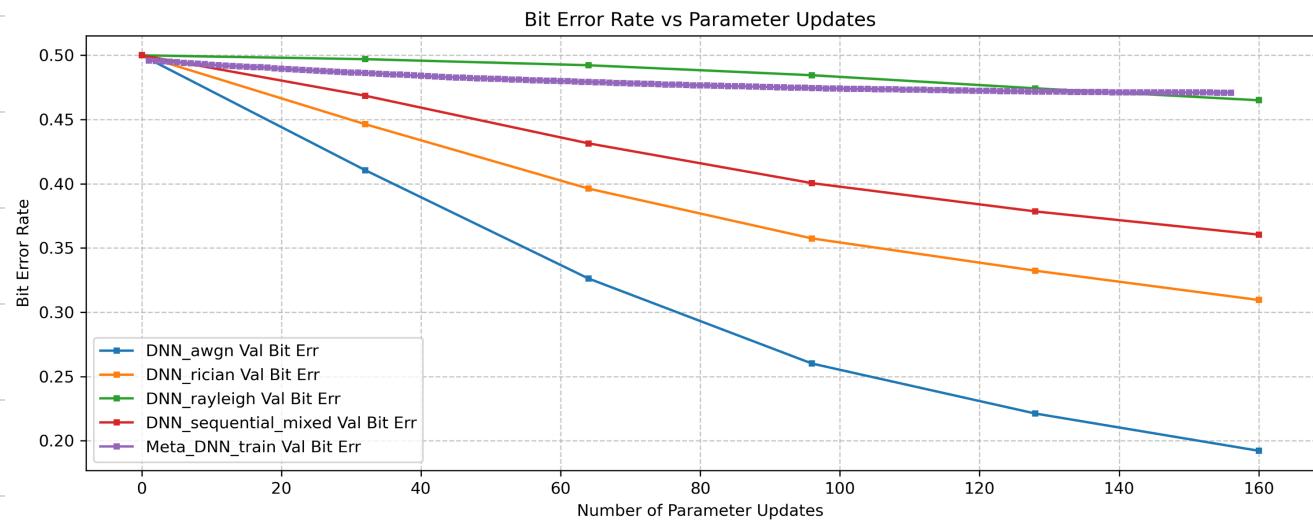
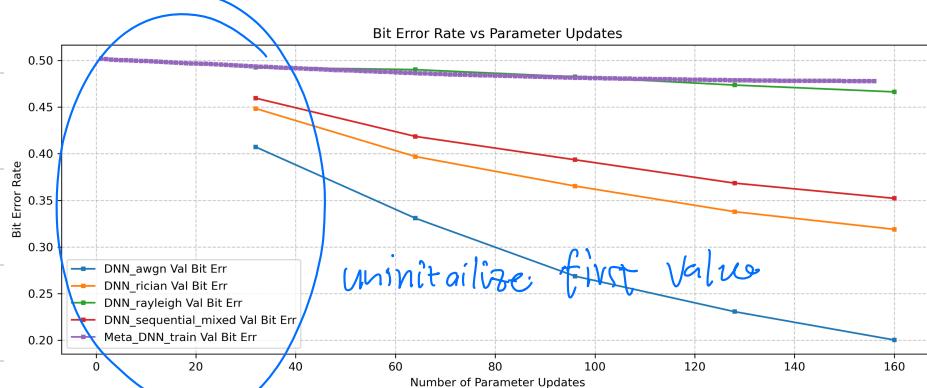
100k total

$$\frac{100k}{128} \times 20$$

mini-batch sampled from task,  
put into inner-update.



trace where it write down  
the Val\_bit by holding Val\_update\_count[SC]



2025/3/5

In fine tune phase, model get polluted

⇒ get cloned before fine tune

Current meta\_lr. Using Linear decay

$$\text{meta\_lr} = \left(1 - \frac{\text{current\_iteration}}{\text{total}}\right) \times \text{meta\_lr}$$

⇒ hyper gradient?

calculate the gradient of  $(\phi' - \phi)$   
and update meta\_lr based on grad.

set  
first value  
 $= 0.5$ .

⇒ Cosine Annealing

$$y = y_{\min} + \frac{1}{2}(y_{\max} - y_{\min}) \left( 1 + \cos\left(\frac{T_{\text{cur}}}{T_{\max}} \pi\right) \right)$$

⇒ use: Cosine Decay Restarts

first\_decay = 100, 450 get converge...

$t_{\text{mul}} = 1.3$ ,  $m_{\text{mul}} = 0.9$ ,  $\alpha = 0.001$

further work: inner\_lr update optimize to avoid overfit or local\_min.

try to decrease train time of Meta-DNN

too update,

$$\text{inner iteration} = \frac{\text{DNN samples}}{\text{minibatch\_size} \times \text{Number of task}} = 27$$

Type

plain

Time

89.43s

1.01 / meta-update

~~@tf.function~~



future wrk

~~parallel~~

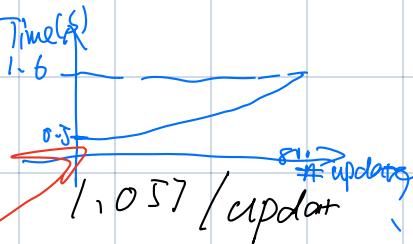
Cosine Restart

~~@tf.function~~

(w inner-update)

667.63s  
1.6 Time(s)

845.78s  
1.057 / update



Tensorflow repeatedly

recompiling



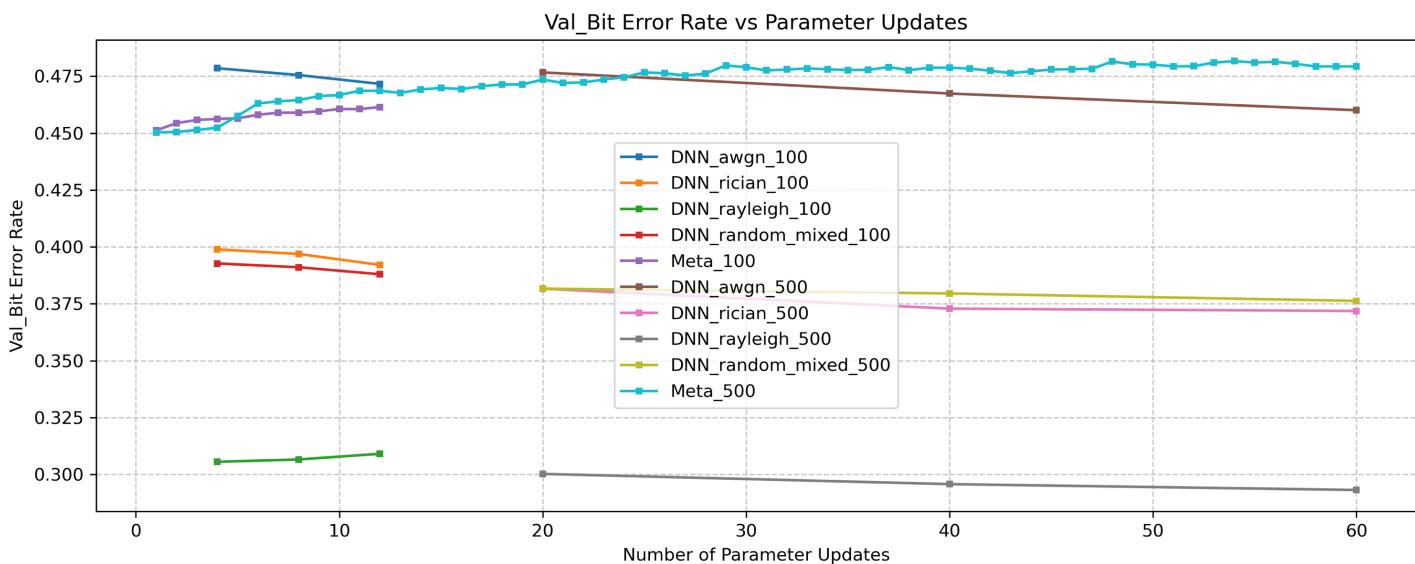
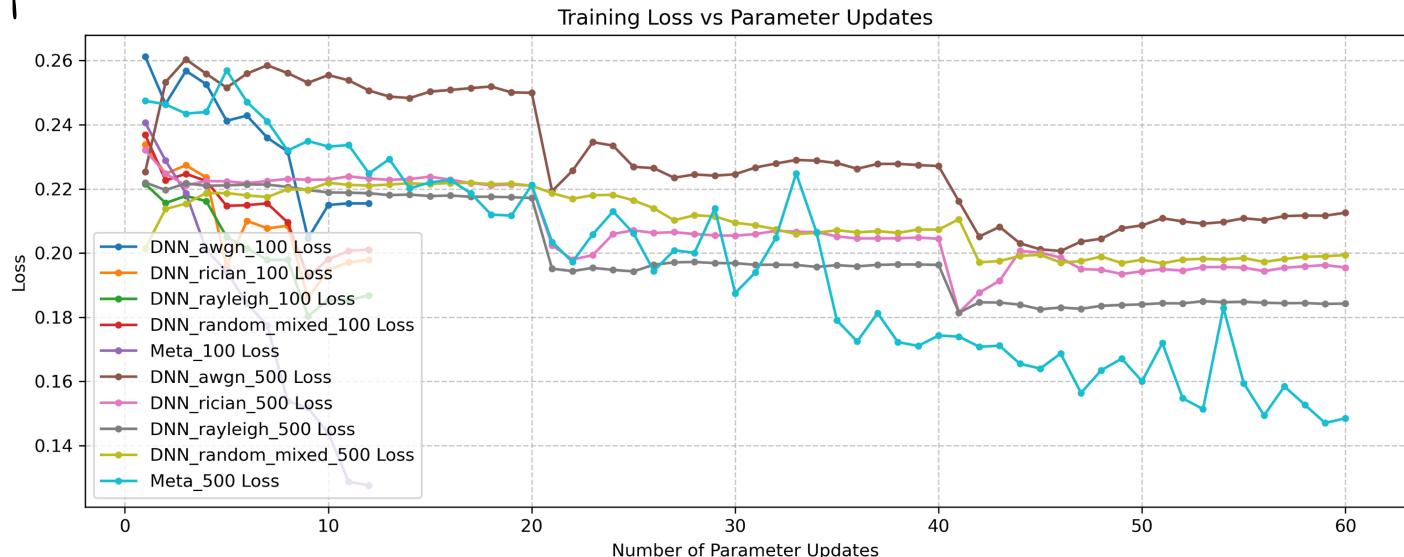
legend issue in

3gpp phase, take average  
or place it outside.

let reduce_retracing = True	version control: DNN function well → 0.1.1
allow tensor to trace graph	update result naming 0.1.2
	add mixed 0.1.3
	Multi BCP 0.2.3
reuse meta_grads	baric metaDNN 0.3.3
to optimize memory management	Train phase 0.4.3
	mini-batch 0.4.4
	training figure (plot-by-update) v. T.4
	Generalize test, convergent, 1.6.5
	Time trace -

Initialize the inner train function with model object and decorated with `tf.function` to utilize graph

update\_time: 800, Meta\_round BER = 0.408828



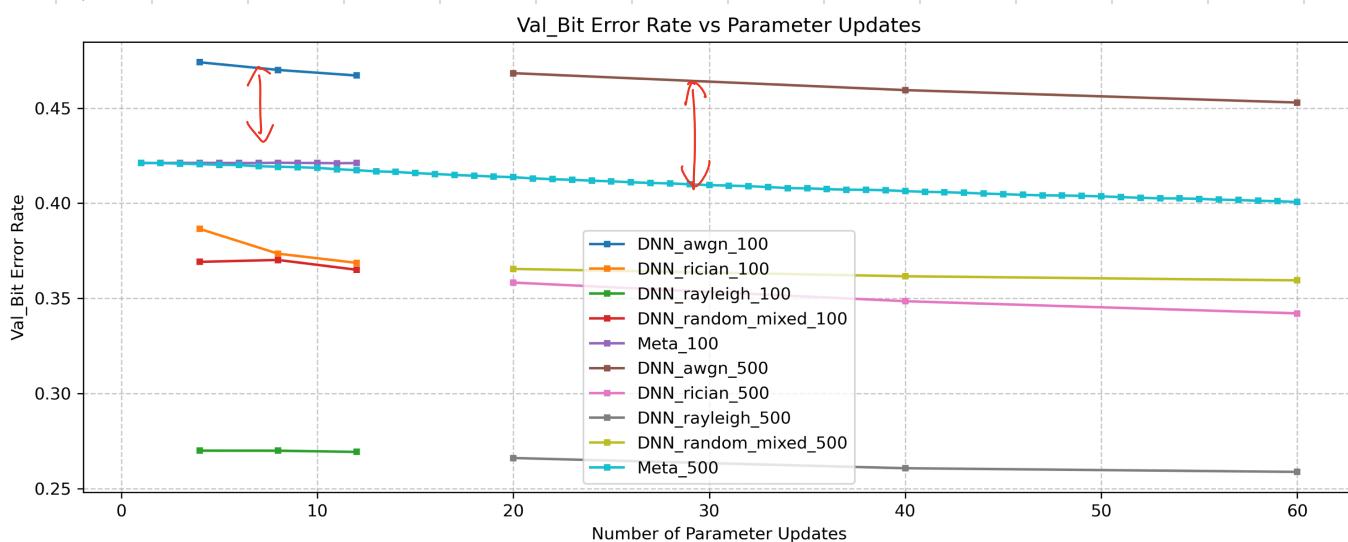
update 1600 → ber trapped at 0.4

remove best\_weight keeper, enumerate update  
→ trapped at 0.4

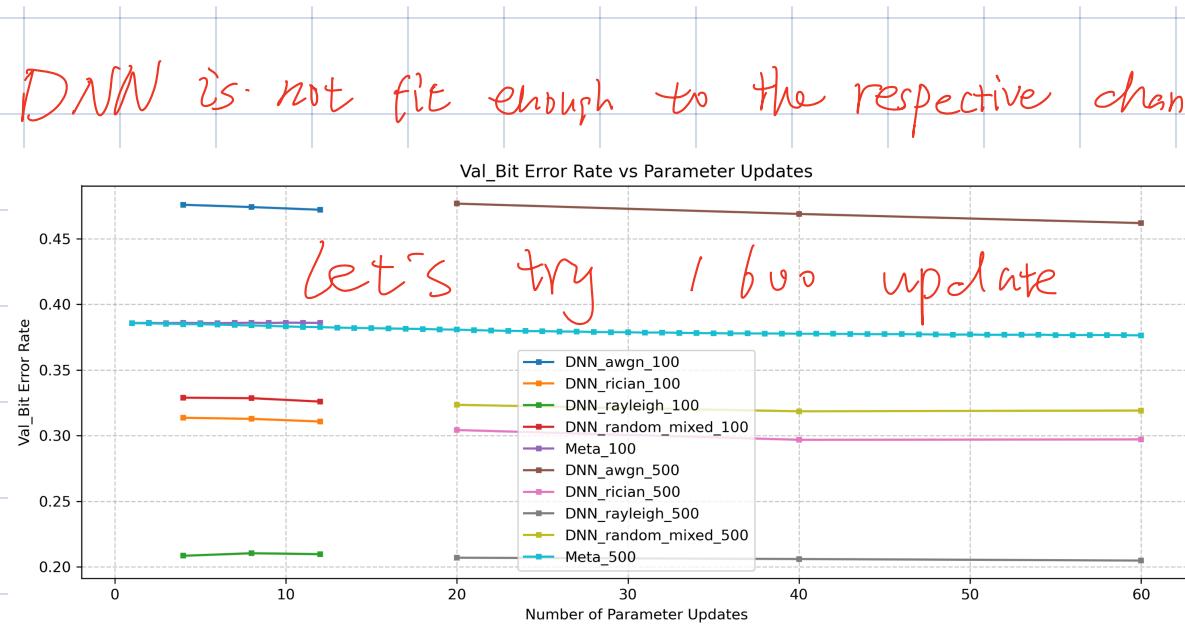
remove tf.function ✓

final keep: Performance tracker, CosineRestarter.

update i 800



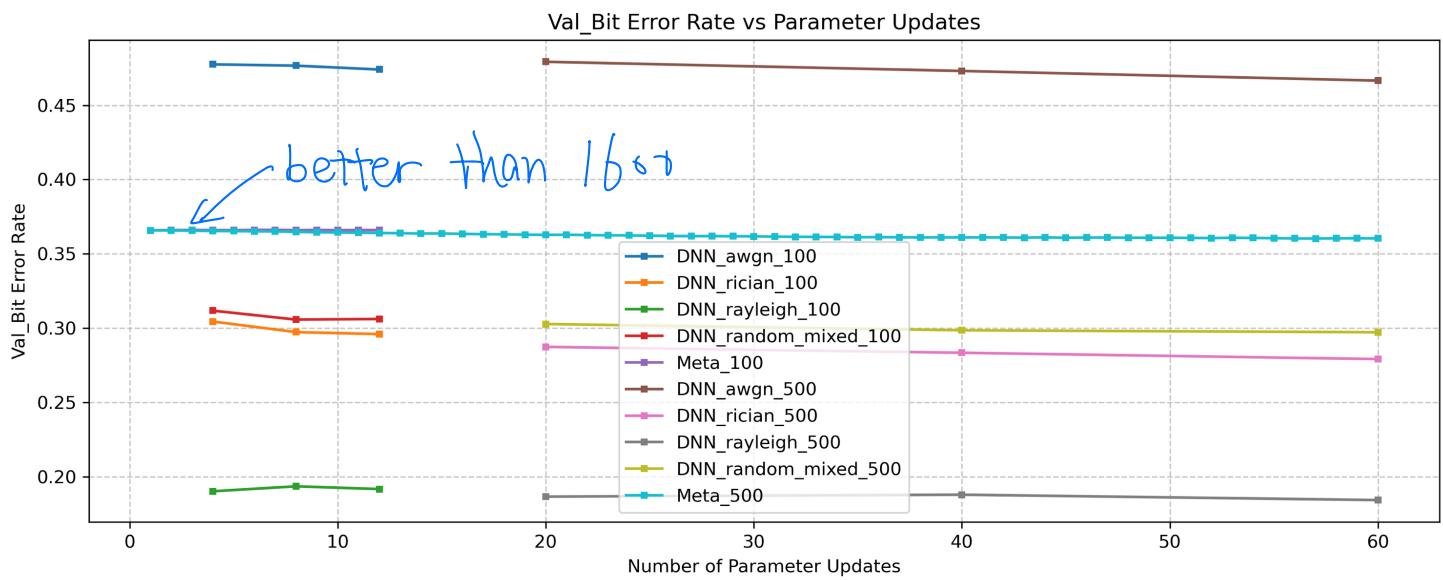
DNN is not fit enough to the respective channels



2024.3.10

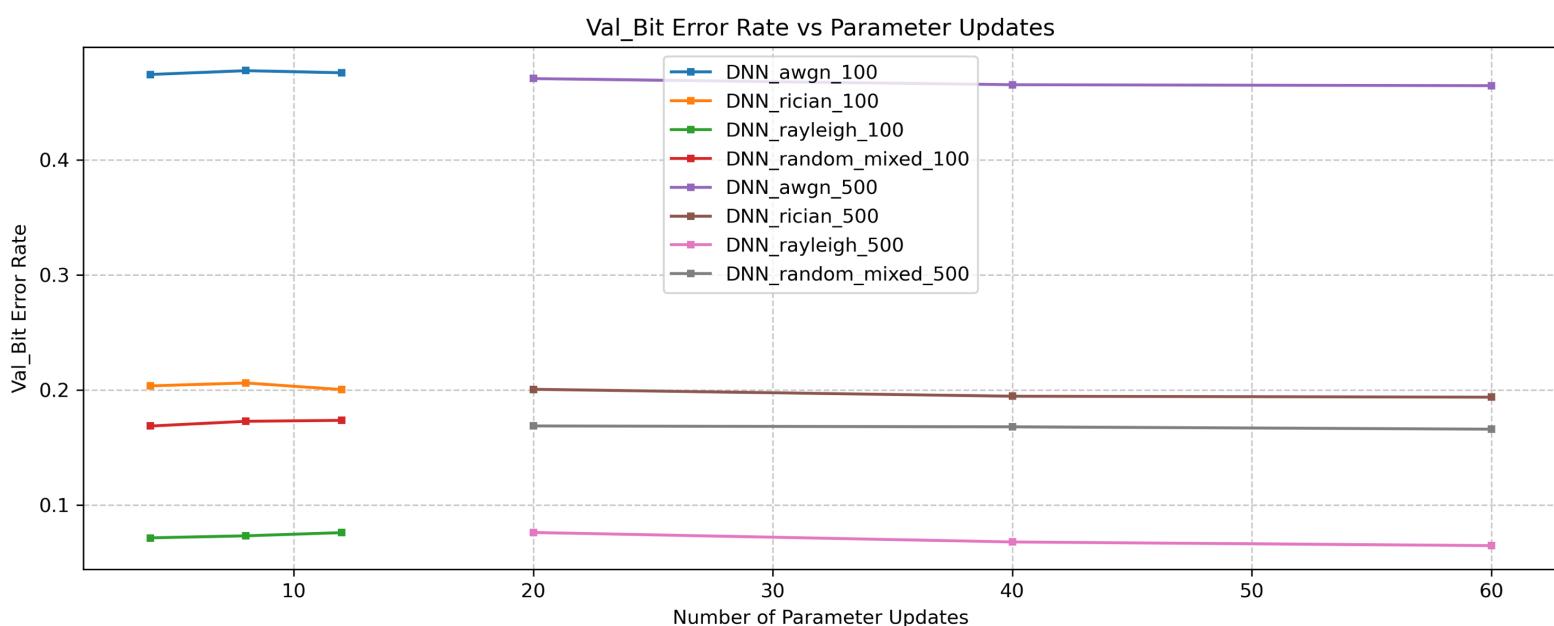
Let's try 2400 updates

What expected is meta-DNN performs better than DNN on 3GPP channel.

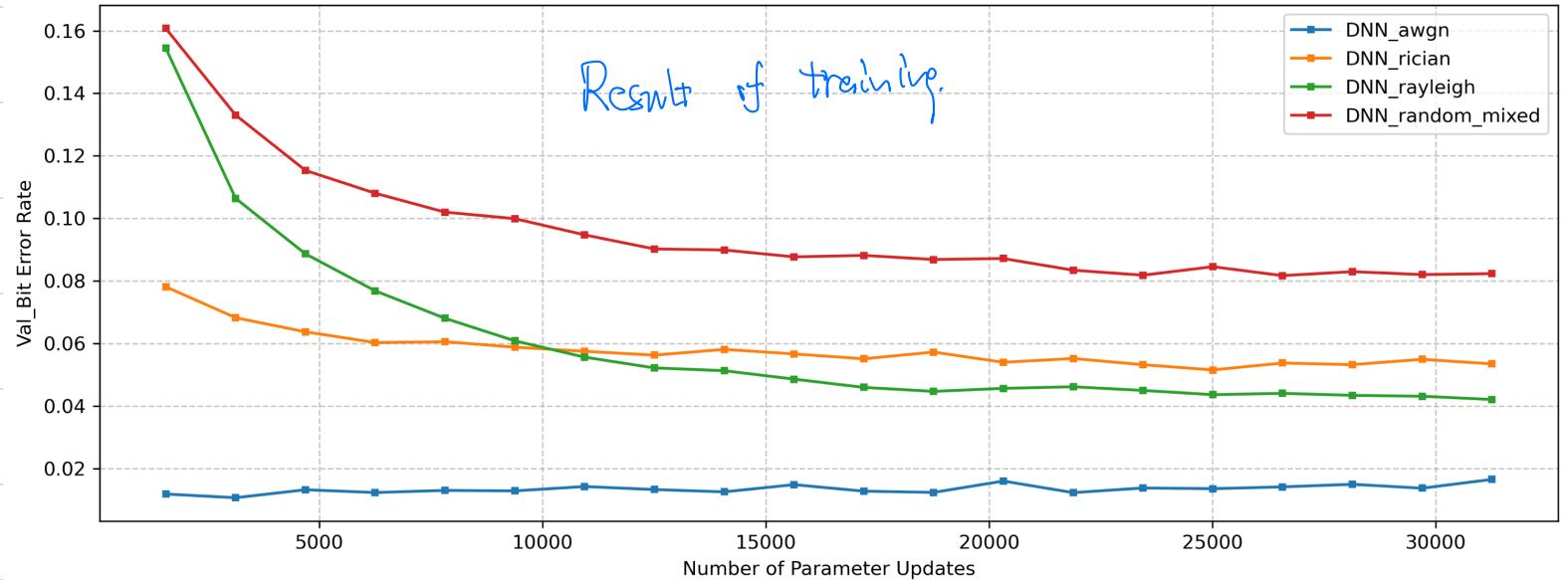


try 10k sample = 31250 update  
(2 epoch, 64 batch size) only DNN

What expected is well-trained DNN will perform worse on 3GPP



Val\_Bit Error Rate vs Parameter Updates



Result of training.

2025/3/11

Supervisor meeting

Train on 3 of them, take rest as target channel

awgn  
rician  
mixed

35PP

fine tune training sample

50  
100  
150  
200

100  
150  
200

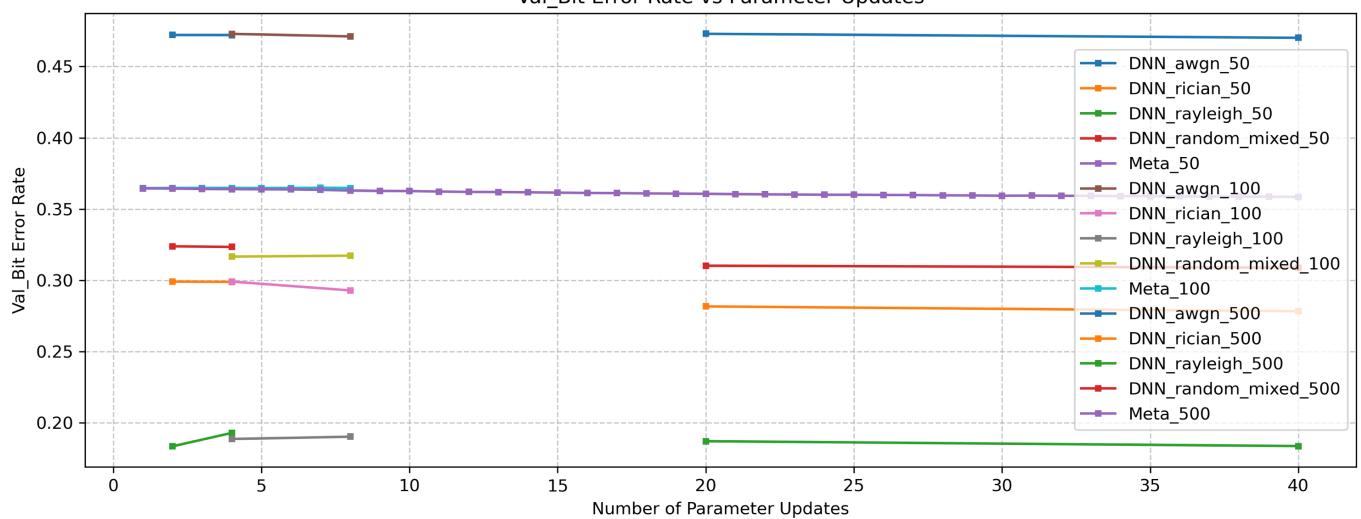
100  
150  
200

100  
150  
200

With 2400 update, add 50 sample set into validation

phase.

Val\_Bit Error Rate vs Parameter Updates



Meta-DNN perform worse than rician, random-mixed and rayleigh generally.

2025/3/15

Review the "On First-order Meta Learning Algorithms"

try cosine Restart applied

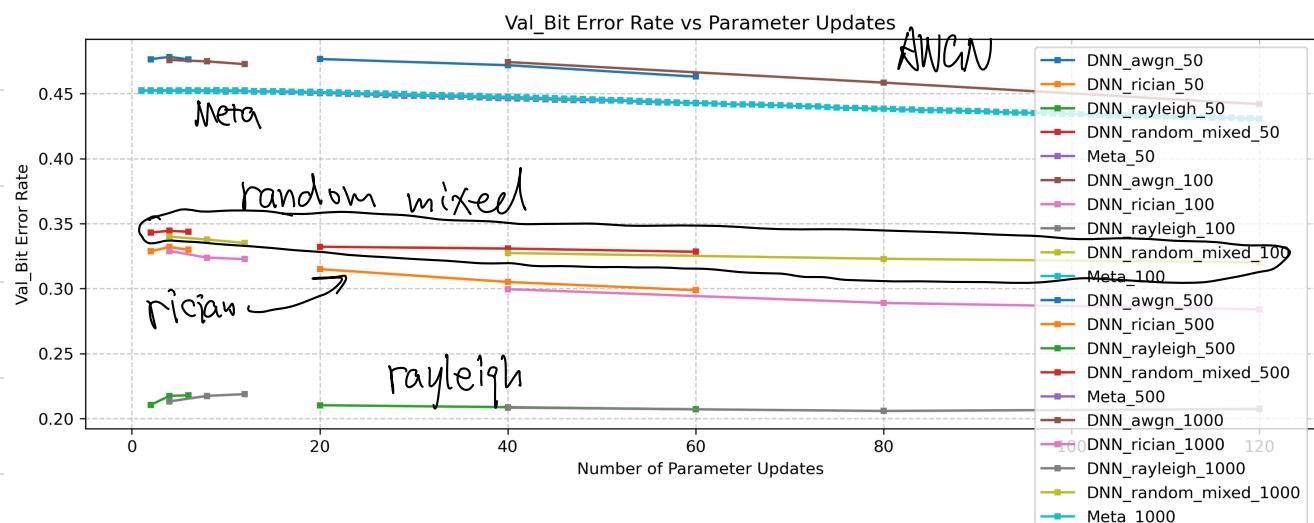
1. apply linear annealing on Meta-update

2. Adjust inner step from {1, 2, 5, 8, 10, 11, 20}

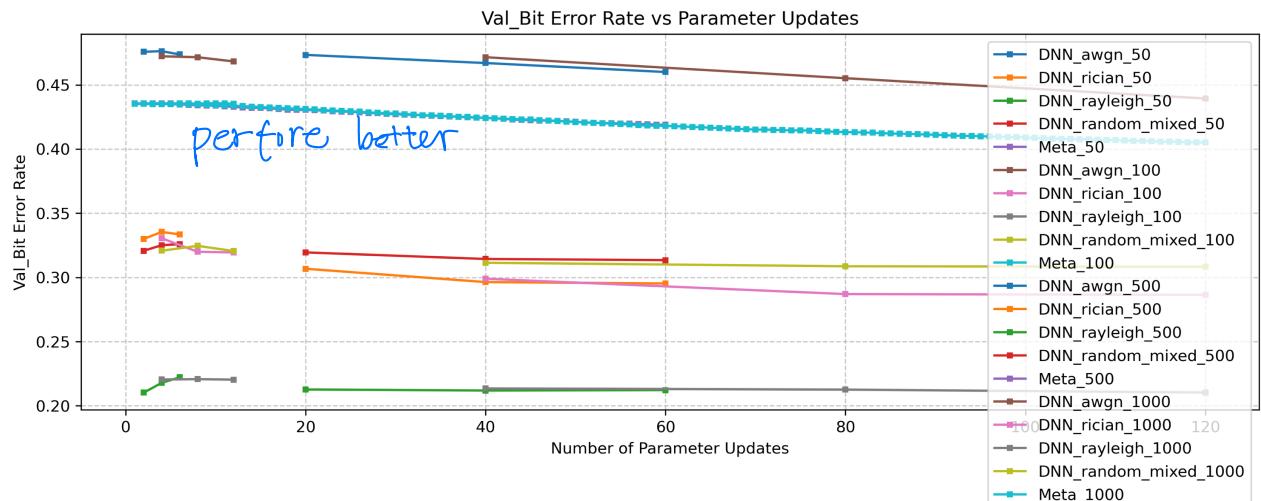
3. mini batch size  $\Rightarrow \{10, 20, 32\}$

task step (inner step) = 8, # of updates = 1600

batch size = 32

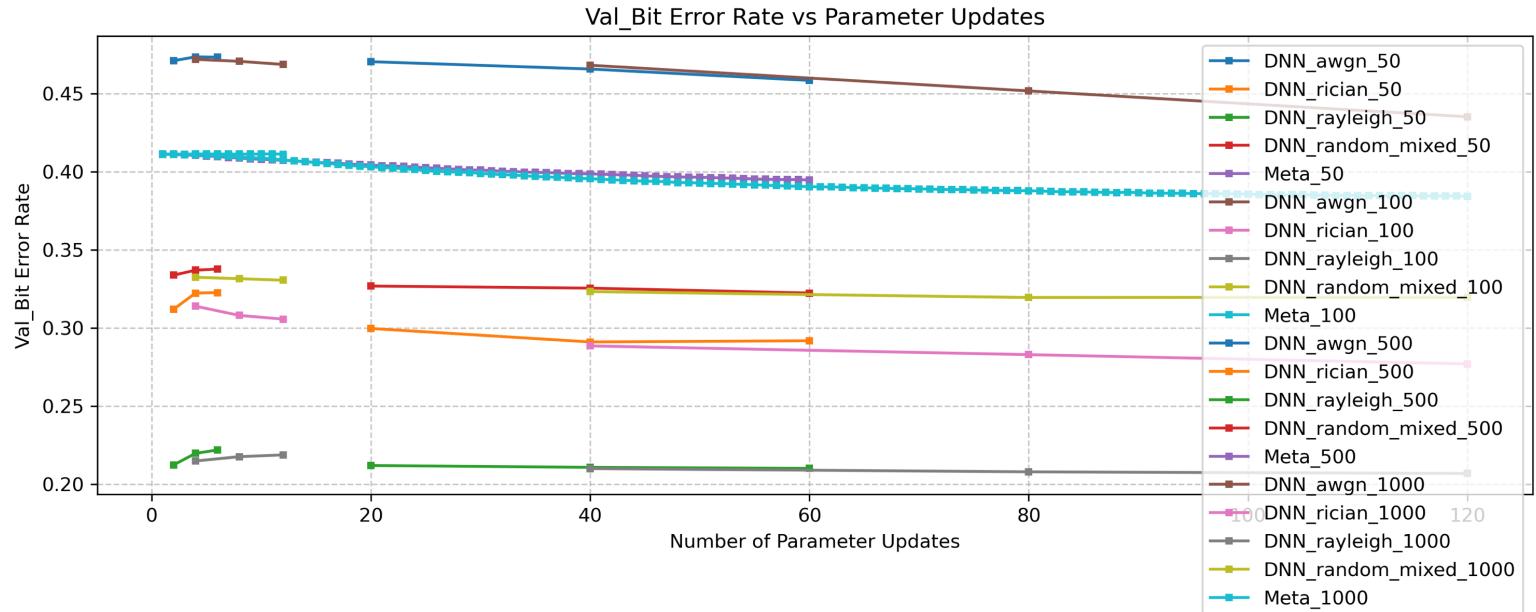


task  
step = 10



task

$$\text{step} = 15$$

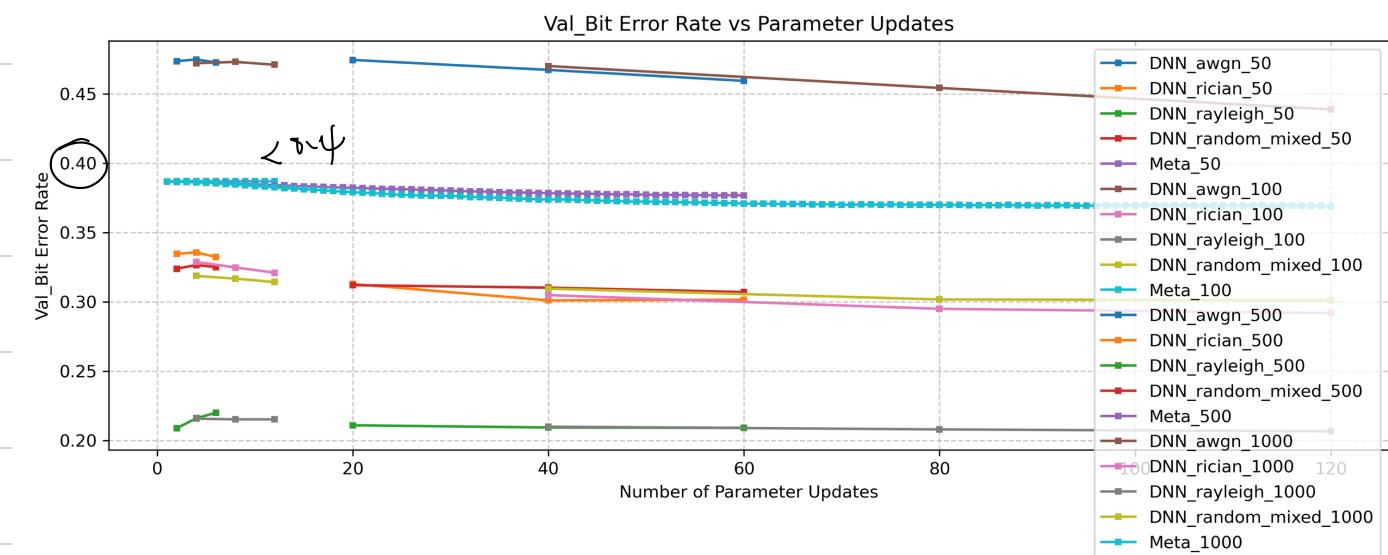


Take Adam in inner update.

-task-step = 15

loss descent very slow.

task - step  $\leq 50$



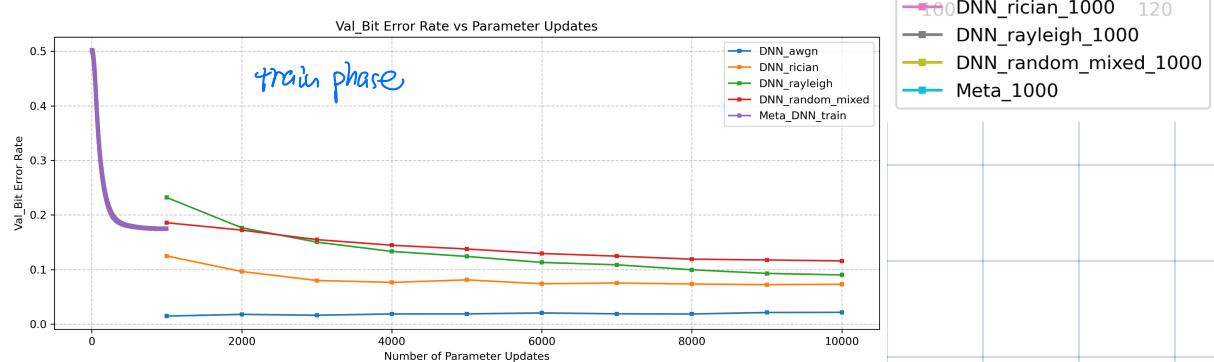
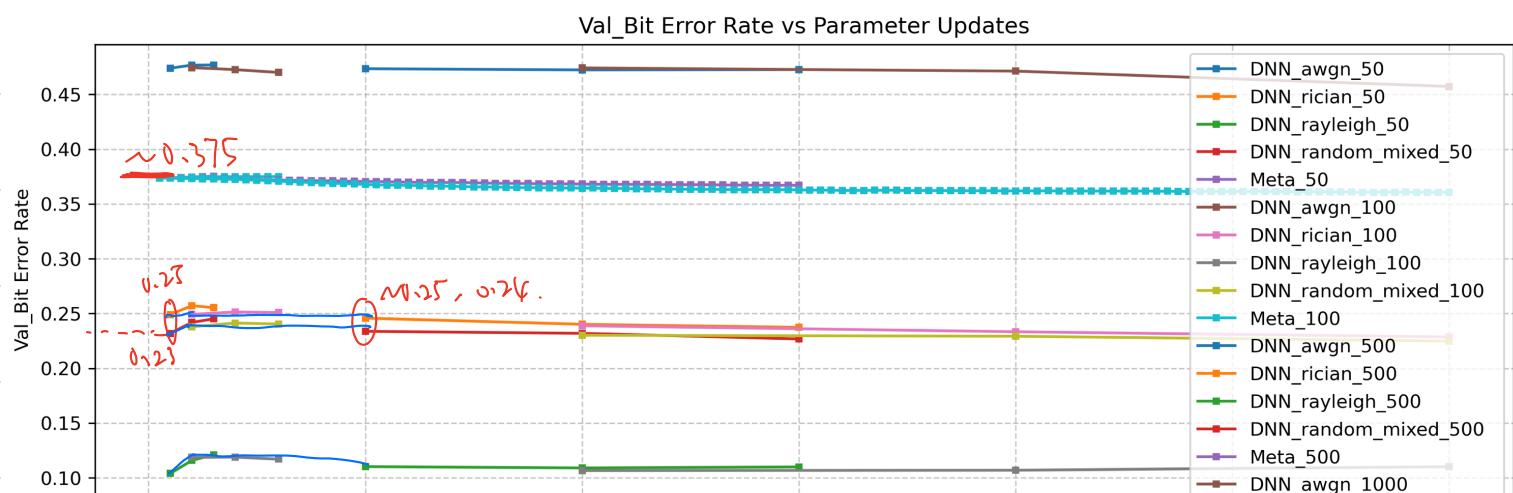
try update time = 10000 (with 32 batch size, 10 epoch)  
innerstep = 50

bit error rate jumping between 0.17~ when 2700~ / 10000  
therefore aborted.

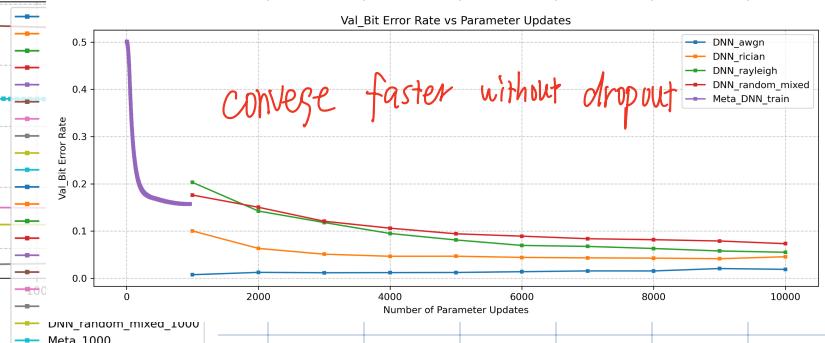
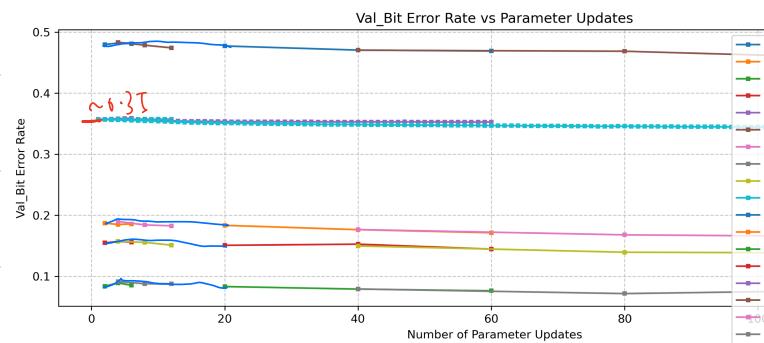
Early stop check added

Patience = 50, if variance < 0.001, patience = patience - 1

Early stop at 996 with BER = 0.174817



What if - remove dropout layer to downgrade the generalization ability of DNN Early stop at 963. with BER = 0.17145



2025/3/17

early stop at 2365 with BER = 0.089622

### Architecture

Hidden: { 256, relu, 512, relu, 128, relu, output: 112, sigmoid }  
 Optimizers = Adam(0.001)  
 Loss = MSE  
 Epoch = 10, batch\_size = 32

### Reptile

Meta-update  
 $\text{meta-lr} = 0.3$ ,

Best para tracer

Warm-up epoch = min(50, Total Iteration // 10) SGD

$$\text{lr} = \text{lr} * \frac{(\text{epoch}+1)}{\text{Warm-up}}$$

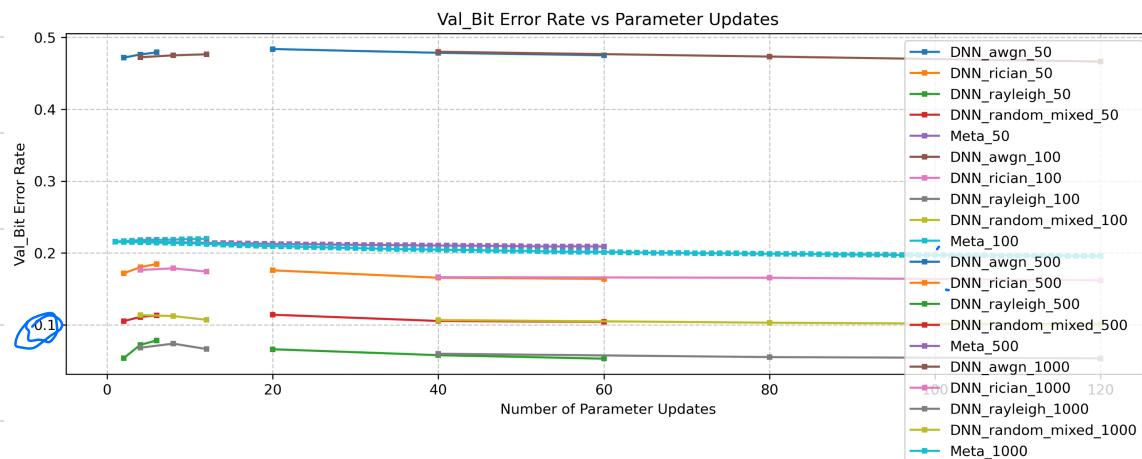
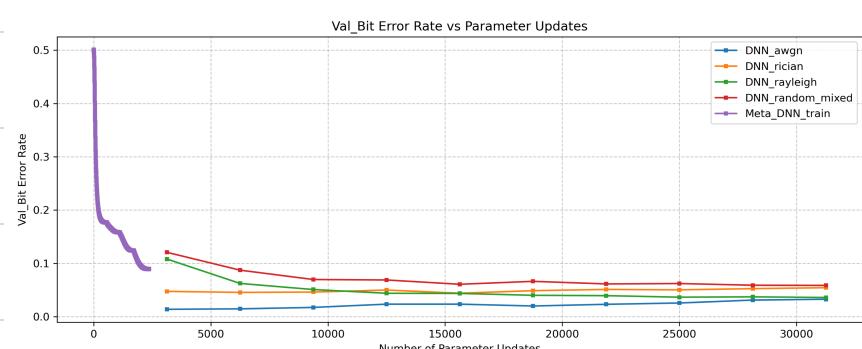
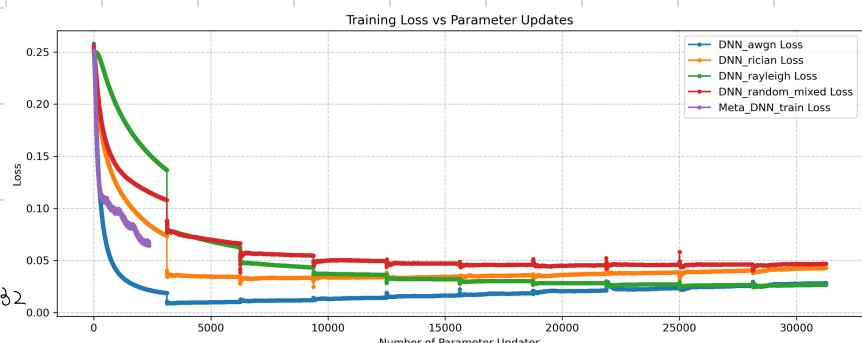
Inner-update  
 $\text{inner-lr} = 0.02$ ,  
 task-step = 50  
 mini-batch-size = 32

MSE

Early stopping & patience = 50,  
 $\text{min-delta} = 0.0005$

Cosine Decay Restart (first-decay-step = 50),

$$\begin{aligned} t\_mul &= 1 / \text{period} & \text{period} &= 1 / \text{X period} \\ m\_mul &= 1 & \text{lr} &= m\_mul \times \text{lr} \\ \alpha &= 0.001 & \text{min-lr} &= \alpha \times \text{initial-lr} \end{aligned}$$



### Second Validation

early stop at 1093 with BER = 0.156662

try remove early stop to fully take advantage of Cosine Decay Restart

2025/3/19

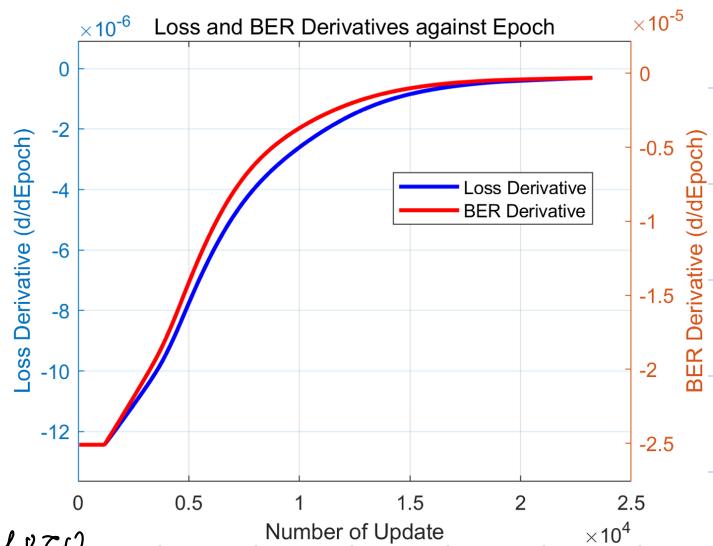
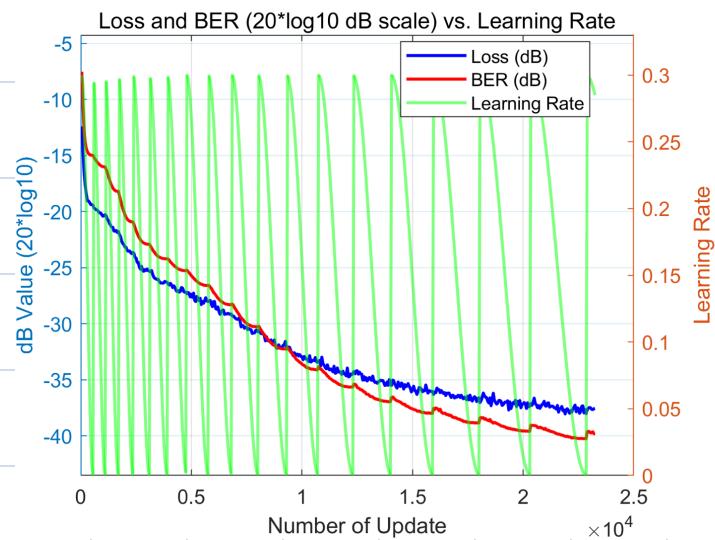
After removing early stop.

Stop by Resource exhausted at 2025/3/25 with

$$BER = 0.01261$$

reveal the relationship between

### Training Process Analysis

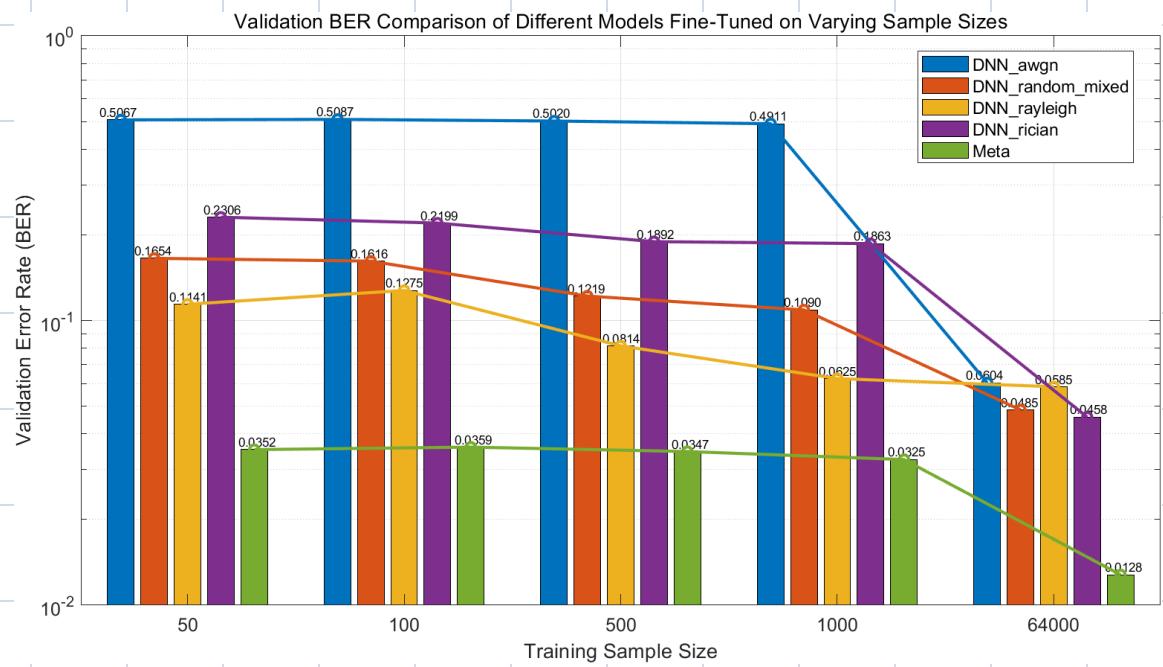


the update progress slow ( $\frac{d BER}{d Epoch} \rightarrow 0$ ) @ update  $\approx 1.5 \times 10^4$

I chose 20K updates for subsequent training.

2025/3/23

Result:



Meta-DNN outperforms than DNNs in limited data case.

$(\delta_0, l_{00}, \Gamma_{00})$