

# Automated Code Transformation for Distributed Training of TensorFlow DL Models

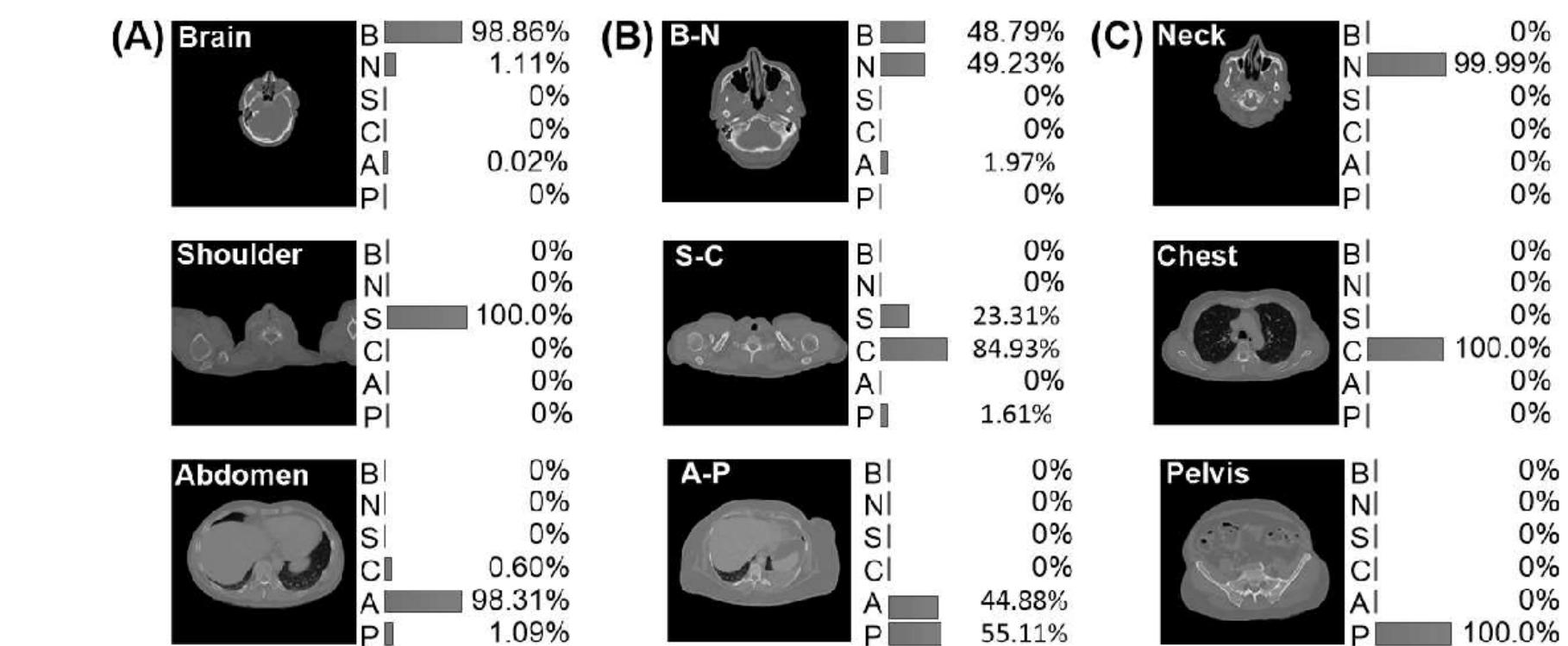
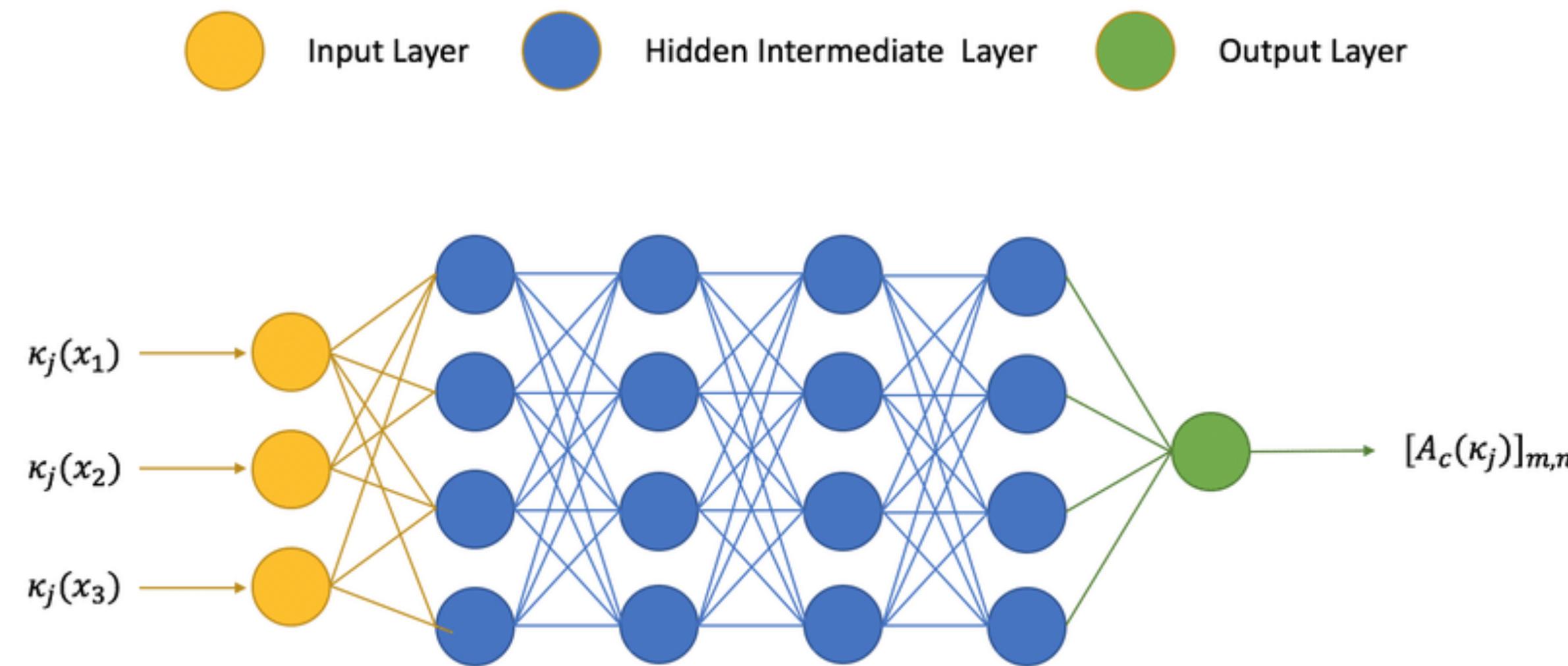
**Yusung Sim<sup>1</sup>, Wonho Shin<sup>1</sup>, Shin Young Ahn<sup>2</sup>, Sungho Lee<sup>3</sup>, Sukyoung Ryu<sup>1</sup>**

<sup>1</sup>School of Computing, KAIST

<sup>2</sup>Artificial Intelligence Research Laboratory / Supercomputing Technology Research Center, ETRI

<sup>3</sup>Department of Computer Science and Engineering, Chungnam National University

# Deep Learning(DL)



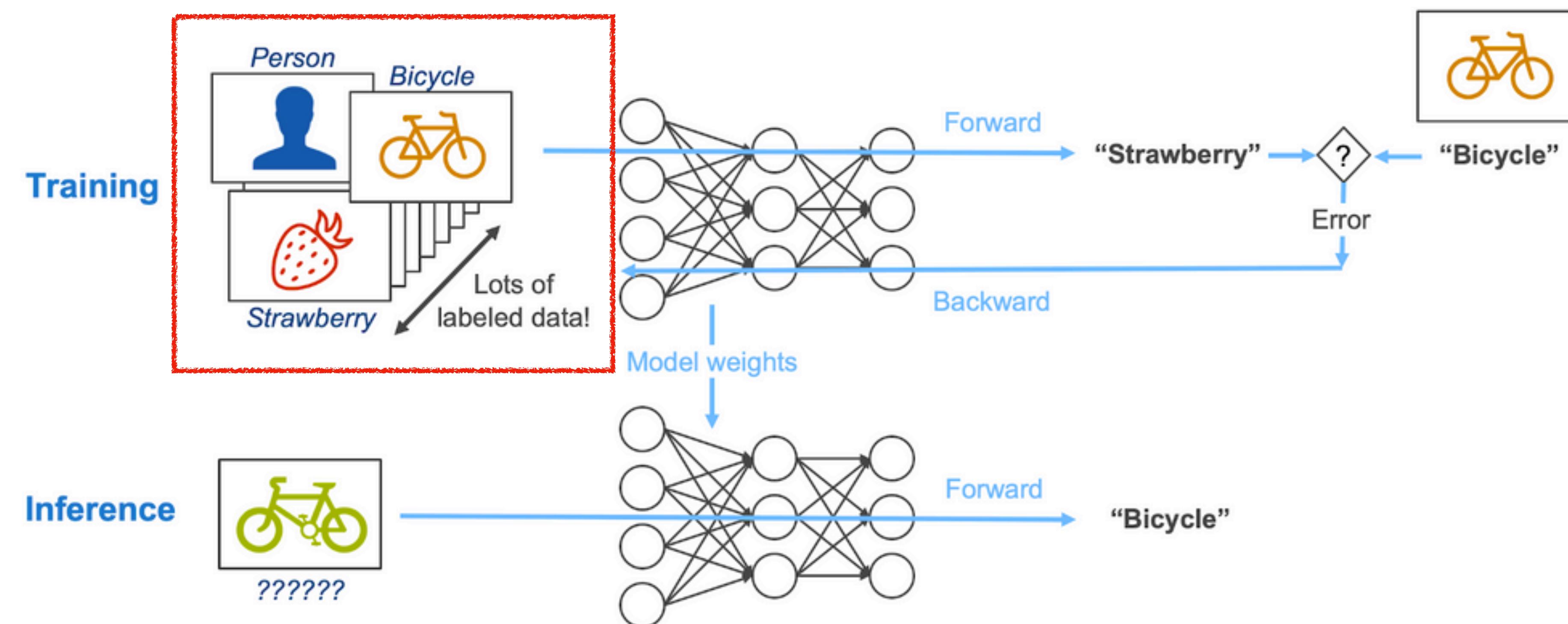
# DL Model Development in Practice



```
1 import tensorflow as tf
2
3 model = tf.keras.models.Sequential([
4     tf.keras.layers.Flatten(input_shape=(28, 28)),
5     tf.keras.layers.Dense(128, activation='relu'),
6     tf.keras.layers.Dropout(0.2),
7     tf.keras.layers.Dense(10)
8 ])
9
10 model.compile(optimizer='adam',
11                 loss=loss_fn,
12                 metrics=['accuracy'])
13 model.fit(x_train, y_train, epochs=5)
```

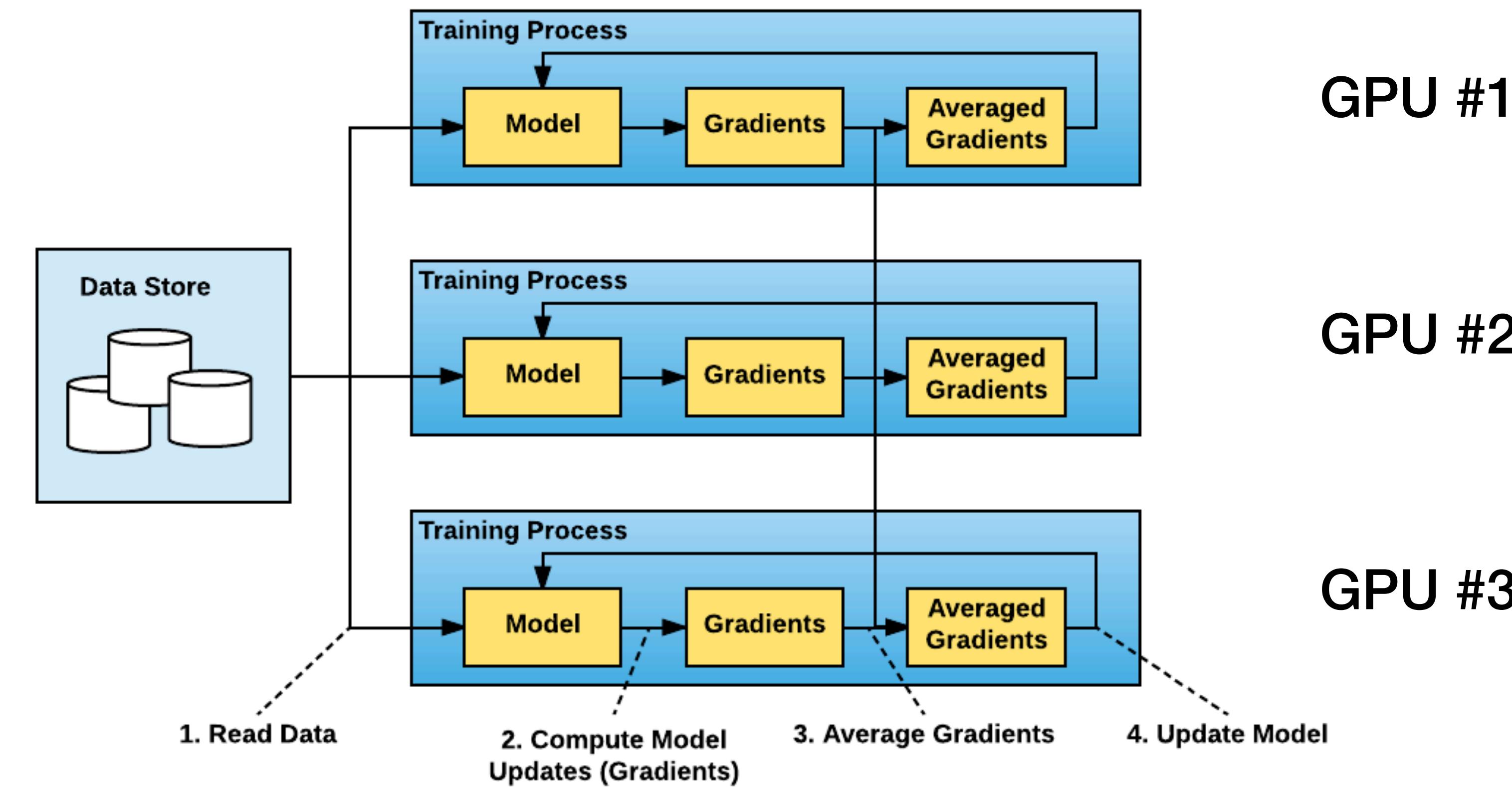
# DL Model Training

- Training is a significant bottleneck in DL development.
  - Repeats training steps for a large number of training data.
  - Requires lots of computation & time.



# Distributed Training

- Technique to reduce the training time by parallelizing the process over multiple GPUs.



# Distributed Training in Practice

## horovod/horovod

Distributed training framework for TensorFlow, Keras, PyTorch, and Apache MXNet.



154  
Contributors

603  
Used by

137  
Discussions

13k  
Stars

2k  
Forks



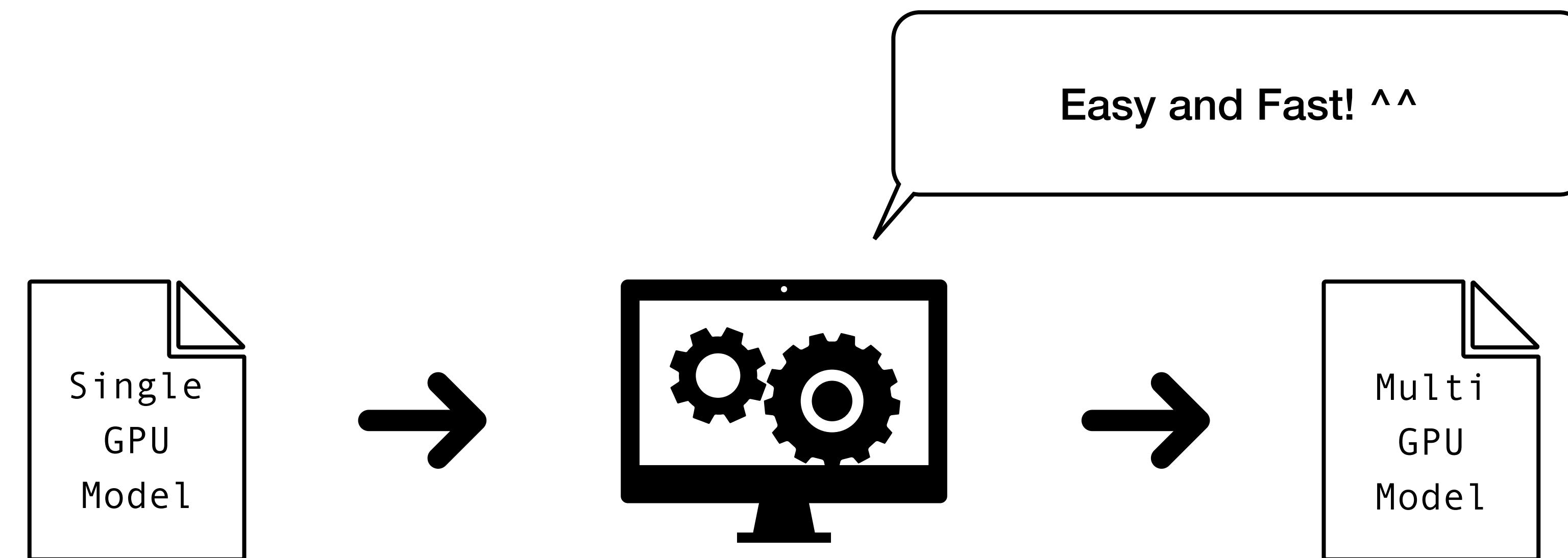
# Challenge: Distributing Existing DL models

- Developers have manually rewritten model codes
- The rewriting is time-consuming & labor-intensive.
- Must understand both original DL models & distributed training library.



# Solution: Automated Code Transformation for Distributed Training

- Automate the rewriting process with the code transformation technique.



# Overview

1. Parse the input model into AST.

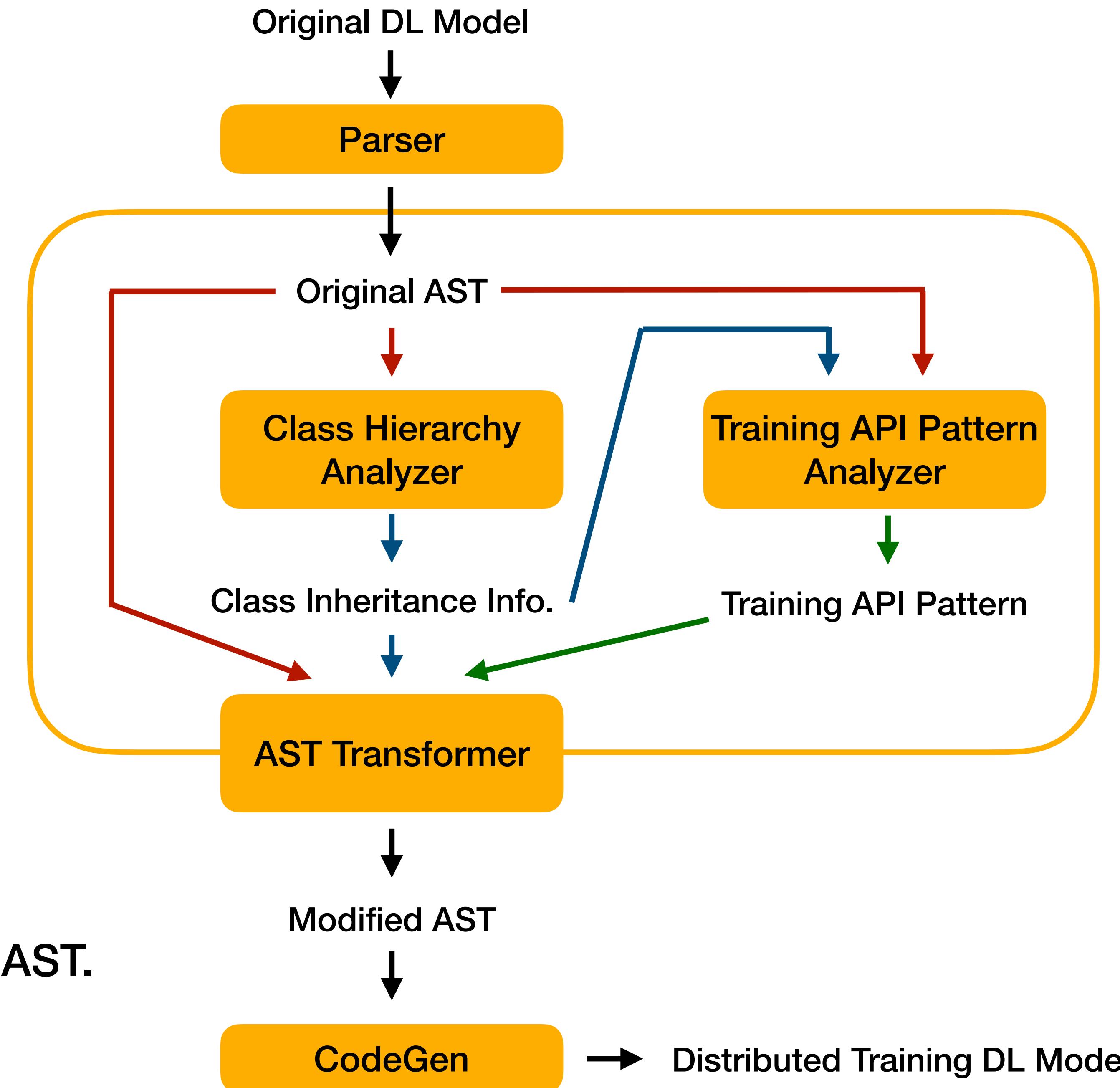
2. Transform the AST.

2.1. Analyze the class hierarchy.

2.2. Analyze the training API pattern.

2.3. Transform the AST,  
using results From 2.1 & 2.2.

3. Generate the distributed model from AST.



# Python Abstract Syntax and Parser

- Full grammar specification in Python Language Reference

<i>module</i>	$::=$	<i>stmt</i>	
<i>stmt</i>	$::=$	<i>def id (args) : stmt*</i>	(FUNDEF)
		<i>class id (expr* keyword*) : stmt*</i>	(CLASSDEF)
		<i>expr = expr</i>	(ASSIGN)
		<i>for expr in expr : stmt*</i>	(FORLOOP)
		<i>while (expr) : stmt*</i>	(WHILELOOP)
		<i>if (expr) : stmt* (else : stmt*)?</i>	(IF)
		<i>with with_item* : stmt*</i>	(WITH)
		<i>import alias*</i>	(IMPORT)
		<i>from i id? import alias*</i>	(IMPORTFROM)
		<i>expr</i>	(EXPRSTMT)
		<i>stmt \n stmt</i>	(SEQUENCE)
<i>expr</i>	$::=$	<i>expr boolop expr</i>	(BOOLOP)
		<i>expr binop expr</i>	(BINARYOP)
		<i>unop expr</i>	(UNARYOP)
		<i>[expr*]</i>	(LIST)
		<i>(expr*)</i>	(TUPLE)
		<i>expr (expr* keyword*)</i>	(CALL)

# Class Hierarchy Analysis

- User-defined classes can inherit TensorFlow library classes.
  - Extend or modify behaviors of the classes or methods.

```
1 from tensorflow import keras
2
3 class ResNet(keras.Model):
4     def __init__(self, block_list):
5         ...
6
7 model = ResNet([2, 2, 2])
8
9 model.fit(x_train, y_train)
```

User-defined class  
for new model

Method inherited from parent class

# Class Hierarchy Analysis

- Correct transformation rule should identify user-defined classes that inherit training-related library classes.

```
1 from tensorflow import keras
2 import horovod.tensorflow.keras as hvd
3
4 class ResNet(keras.Model):
5     def __init__(self, block_list):
6         ...
7
8     model = ResNet([2, 2, 2])
9
10    model.fit(x_train, y_train,
11               callbacks=[hvd.callbacks.BroadcastGlobalVariablesCallback(0)])
```

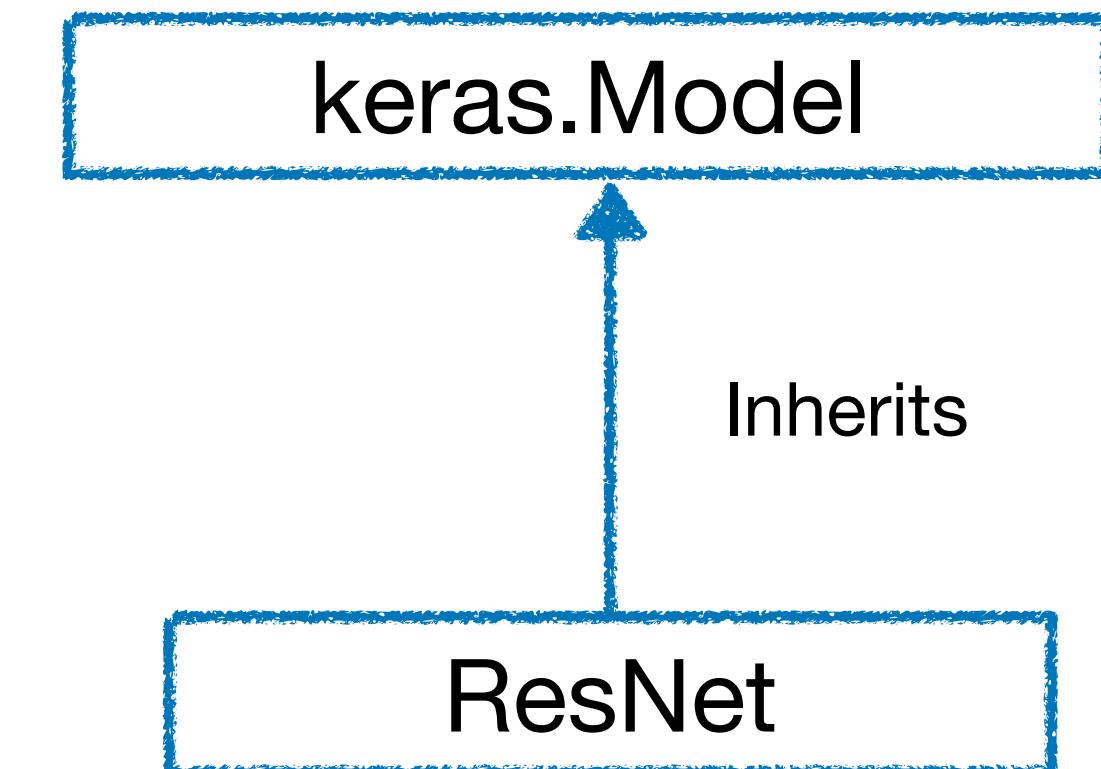
Must know that this instance  
Inherits 'fit' from 'keras.Model'

Adds argument for the 'fit' method

# Class Hierarchy Analysis

- Pre-analyze class inheritance relation between TensorFlow library classes and user-defined classes.
- The class hierarchy analysis result is sent to later stages of the transformation.

```
1 from tensorflow import keras
2 import horovod.tensorflow.keras as hvd
3
4 class ResNet(keras.Model):
5     def __init__(self, block_list):
6         ...
7
8 model = ResNet([2, 2, 2])
9
10 model.fit(x_train, y_train,
11             callbacks=[hvd.callbacks.BroadcastGlobalVariablesCallback(0)])
```



# Training API Pattern Analysis

- TensorFlow provides different APIs to define the training.
- Different usage of training APIs requires different transformation rules.

```
1 for x, y in train_data.take(training_steps):  
2     with tf.GradientTape() as tape:  
3         pred = model(x, is_training = True)  
4         loss = loss_compute(y, pred)  
5  
6     trainable_vars = model.trainable_variables  
7     gradients = tape.gradient(loss, trainable_vars)  
8     pairs = zip(gradients, trainable_vars)  
9     optimizer.apply_gradients(pairs)
```

Repeating training steps with ‘for’ loop

```
1 model.compile(  
2     optimizer = optimizer,  
3     loss = loss_compute)  
4 model.fit(train_data.take(training_steps))
```

Calling ‘fit’ method to  
start the automatic training process

# Training API Pattern Analysis

- Define four training API patterns categorizing TensorFlow DL models.
- Define a transformation rule for each API pattern.

TF version	Pattern name	Explanation
1.x	Session	Low-level training API using <code>Session</code> instance
1.x	MonitoredSession	Low-level training API using <code>MonitoredSession</code> instance
2.x	GradientTape	Low-level training API using <code>GradientTape</code> instance
2.x	Keras	High-level training API using <code>fit</code> method of <code>keras.models.Model</code> instance

**FIGURE 7** Training API patterns

# Session Pattern

```
1 with tf.Session() as sess:  
2     for step in xrange(int(num_epochs * train_size) // BATCH_SIZE):  
3         sess.run(optimizer, feed_dict=feed_dict)
```

- Use a Session instance to manually repeat training steps.
- Search a with statement creating a Session instance.

# MonitoredSession Pattern

```
1 with tf.train.MonitoredTrainingSession(hooks=hooks) as mon_sess:  
2     while not mon_sess.should_stop():  
3         mon_sess.run(train_op, feed_dict=feed_dict)
```

- Use a MonitoredSession instance to manually repeat training steps.
  - Similar to the Session pattern.
  - Can add hooks: actions that automatically execute for each step.
  - Search a with statement that creates a MonitoredSession instance.

# GradientTape Pattern

```
1 for x, y in train_data.take(training_steps):
2     with tf.GradientTape() as g:
3         pred = conv_net(x)
4         loss = cross_entropy(pred, y)
5
6     trainable_variables = list(weights.values()) + list(biases.values())
7     gradients = g.gradient(loss, trainable_variables)
8     optimizer.apply_gradients(zip(gradients, trainable_variables))
```

- Use a GradientTape instance to manually repeat training steps.
  - Record operations inside with statement, then automatically compute gradients to optimize model parameters with apply\_gradients method.
  - Search a with statement that creates a GradientTape instance.

# Keras Pattern

```
1 # model definition
2 class ResNet(tf.keras.models.Model):
3     def __init__(self, params):
4         ...
5
6 model = ResNet([2, 2, 2], num_classes)
7 model.fit(x_train, y_train_ohe, batch_size=batch_size, epochs=epochs,
```

- Define a model with Keras library APIs, and use the `fit` method for automatic training.
- Search for the `Model.fit` method call.
  - Utilize class inheritance information to identify user-defined model classes.

# AST Transformation Rule

- Manually inspected Horovod documentation and code examples.
- Define four transformation rules, each for training API patterns.

$$\text{trans}_S : \text{Stmt} \rightarrow \Sigma \rightarrow (\text{Stmt list} \times \Sigma)$$

$\text{trans}_S \llbracket id_r = expr_1 (expr_{11} \dots expr_{1n} (id_1 = )? expr_{21} \dots (id_k = )? expr_{2k}) \rrbracket(\sigma) =$   
**IF**  $expr_1 <: \text{tensorflow.keras.optimizers.Optimizer}$  **THEN** Pattern matching  
for function argument  
**IF**  $id_i = \text{learning\_rate}$  **WHEN**  $1 \leq i \leq k$  **THEN**  
 $([id_r = expr_1 (expr_{11} \dots expr_{1n} (id_1 = )? expr_{21} \dots id_i = expr_{2i} * \text{hvd.size}()$   
 $\dots (id_k = )? expr_{2k}) (\#type:s)?$   
 $id_r = \text{hvd.DistributedOptimizer}(id_r)], \sigma["\text{optimizer"} \mapsto id_r])$

Part of formalization of code transformation rule.

# Transformation rule for Session pattern

```
1 optimizer = tf.train.MomentumOptimizer(  
2     learning_rate = 0.01)  
3  
4 with tf.Session() as sess:  
5     for step in range(num_epochs):  
6         sess.run(optimizer, feed_dict)
```

Original Session pattern training code

```
1 optimizer = tf.train.MomentumOptimizer(  
2     learning_rate = 0.01 * hvd.size())  
3 optimizer = hvd.DistributedOptimizer(optimizer)  
4  
5 with tf.Session() as sess:  
6     for step in range(num_epochs):  
7         sess.run(optimizer, feed_dict)
```

Transformed Session pattern training code

- Modify the Optimizer instance.
  - Find a statement that assigns the result of Optimizer class constructors.
  - Scale the `learning_rate` argument by `hvd.size()`.
  - Wrap the Optimizer instance with `DistributeOptimizer()`.

# Transformation rule for MonitoredSession pattern

```
1 with MonitoredTrainingSession(hooks=hooks) as mon_sess:  
2     while not mon_sess.should_stop():  
3         mon_sess.run()
```

Original MonitoredSession pattern training code

```
1 with MonitoredTrainingSession(  
2     hooks=hooks.append(hvd.BroadcastGlobalVariablesHook(0))) as mon_sess:  
3     while not mon_sess.should_stop():  
4         mon_sess.run()
```

Transformed MonitoredSession pattern training code

- Add the hook that broadcasts global variables over multiple processes.
  - Find a with statement that creates a MonitoredSession instance.
  - Find (or add) the hooks keyword argument.
  - Append the BroadcastGlobalVariablesHook to the argument.

# Transformation rule for GradientTape pattern

```
1 import tensorflow as tf
2
3 with tf.GradientTape() as tape:
4     probs = model(images)
5     loss_value = loss(labels, probs)
6
7 grads = tape.gradient(loss_value, model.trainable_variables)
8 opt.apply_gradients(zip(grads, model.trainable_variables))
```

Original GradientTape pattern training code

```
1 import tensorflow as tf
2 import horovod.tensorflow as hvd
3 hvd.broadcast_done = False
4
5 with tf.GradientTape() as tape:
6     probs = model(images)
7     loss_value = loss(labels, probs)
8     tape = hvd.DistributedGradientTape(tape)
9     grads = tape.gradient(loss_value, model.trainable_variables)
10    id_new = zip(grads, model.trainable_variables)
11    opt.apply_gradients(id_new)
12
13 global hvd.broadcast_done
14 if not hvd.broadcast_done:
15     hvd.broadcast_variables([x[1] for x in id_new], root_rank=0,
16 )     hvd.broadcast_variables(opt.variables(), root_rank=0, )
17     hvd.broadcast_done = True
```

Transformed GradientTape pattern training code

- Modify the GradientTape instance.
  - Find a with statement that creates a GradientTape instance.
  - Wrap the instance with DistributedGradientTape( ).

# Transformation rule for GradientTape pattern

```
1 import tensorflow as tf
2
3 with tf.GradientTape() as tape:
4     probs = model(images)
5     loss_value = loss(labels, probs)
6
7 grads = tape.gradient(loss_value, model.trainable_variables)
8 opt.apply_gradients(zip(grads, model.trainable_variables))
```

Original GradientTape pattern training code

```
1 import tensorflow as tf
2 import horovod.tensorflow as hvd
3 hvd.broadcast_done = False
4
5 with tf.GradientTape() as tape:
6     probs = model(images)
7     loss_value = loss(labels, probs)
8 tape = hvd.DistributedGradientTape(tape)
9 grads = tape.gradient(loss_value, model.trainable_variables)
10 id_new = zip(grads, model.trainable_variables)
11 opt.apply_gradients(id_new)
12
13 global hvd.broadcast_done
14 if not hvd.broadcast_done:
15     hvd.broadcast_variables([x[1] for x in id_new], root_rank=0,
16 )     hvd.broadcast_variables(opt.variables(), root_rank=0, )
17     hvd.broadcast_done = True
```

Transformed GradientTape pattern training code

- Add variable broadcasting code after `apply_gradients` method call.
  - Extract the model parameter list from the `apply_gradients` method argument.
  - Add a fixed set of codes after the call.

# Transformation rule for Keras pattern

```
1 class ResNet(keras.Model):  
2     def __init__(self, block_list):  
3         ...  
4  
5 optim = keras.optimizers.Adam(0.001)  
6  
7 model = ResNet([2, 2, 2])  
8  
9 model.fit(x_train, y_train)
```

Original Keras pattern training code

```
1 class ResNet(keras.Model):  
2     def __init__(self, block_list):  
3         ...  
4  
5 optim = keras.optimizers.Adam(0.001 * hvd.size(), )  
6 optim = hvd.DistributedOptimizer(optim, )  
7  
8 model = ResNet([2, 2, 2])  
9  
10 model.fit(x_train, y_train,  
11             callbacks=[hvd.callbacks.BroadcastGlobalVariablesCallback(0)])
```

Transformed Keras pattern training code

- Modify the Optimizer instance.
  - Find a statement that assigns the result of Optimizer class constructors.
  - Scale the learning\_rate argument by hvd.size( ).
  - Wrap the Optimizer instance with DistributeOptimizer( ).

# Transformation rule for Keras pattern

```
1 class ResNet(keras.Model):  
2     def __init__(self, block_list):  
3         ...  
4  
5 optim = keras.optimizers.Adam(0.001)  
6  
7 model = ResNet([2, 2, 2])  
8  
9 model.fit(x_train, y_train)
```

```
1 class ResNet(keras.Model):  
2     def __init__(self, block_list):  
3         ...  
4  
5 optim = keras.optimizers.Adam(0.001 * hvd.size(), )  
6 optim = hvd.DistributedOptimizer(optim, )  
7  
8 model = ResNet([2, 2, 2])  
9  
10 model.fit(x_train, y_train,  
11             callbacks=[hvd.callbacks.BroadcastGlobalVariablesCallback(0)])
```

- Add the hook that broadcasts global variables over multiple processes.
  - Find the `model.fit` method call.
  - Modify (or add) the `callbacks` keyword argument by adding the `BroadcastGlobalVariablesCallback` hook to the list.

# Evaluation: Transformation

- Experiment on 16 open-source DL models.
- Transformation succeeds on 15/16 models.

Model Name	API Pattern	Transformation Success
LSTM-MNIST	GradientTape	o
SimpleCNN-GradientTape-1	GradientTape	x
SimpleCNN-GradientTape-2	GradientTape	o
SimpleCNN-MonitoredSession	MonitoredSession	o
SimpleCNN-Session	Session	o
VGG-CIFAR10	Keras	o
Play-with-MNIST	GradientTape	o
Linear-Regression	GradientTape	o
Fashion-MNIST	Keras	o
CIFAR10-VGG16	GradientTape	o
Inception-Network	GradientTape	o
RNN-Sentiment-Analysis	Keras	o
Stacked-LSTM-ColorBot	GradientTape	o
Auto-Encoder	GradientTape	o
Variational-Auto-Encoder	GradientTape	o
DCGAN	GradientTape	o

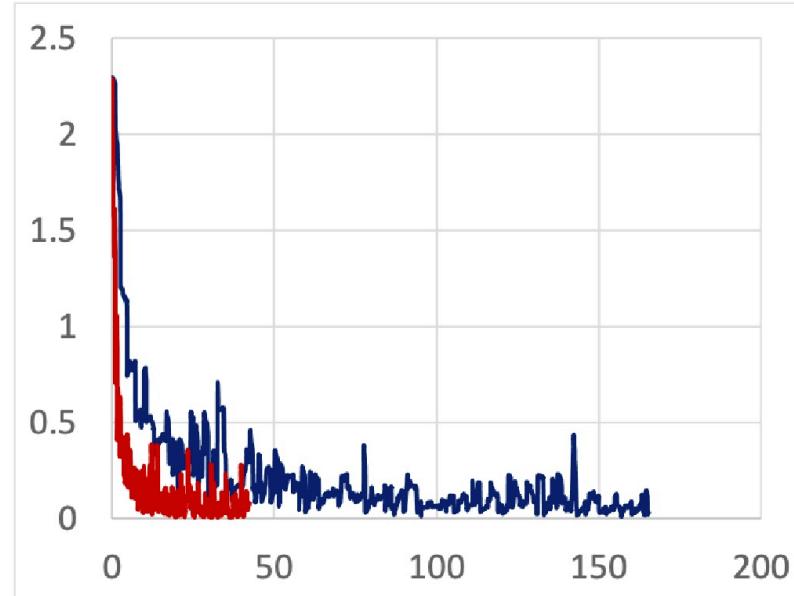
**FIGURE 24** Transformation experiment Result

# Evaluation: Training

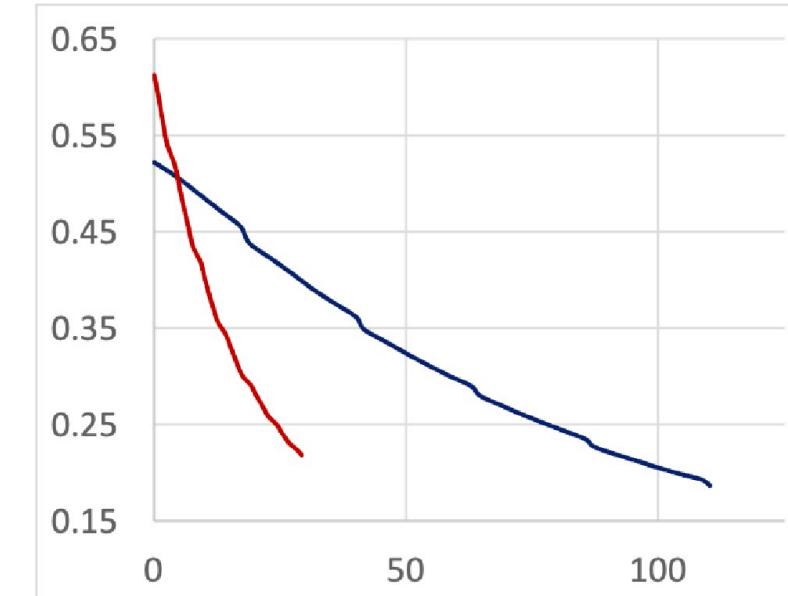
Speed-up: 7/13

No speed-up or Slower: 3/13

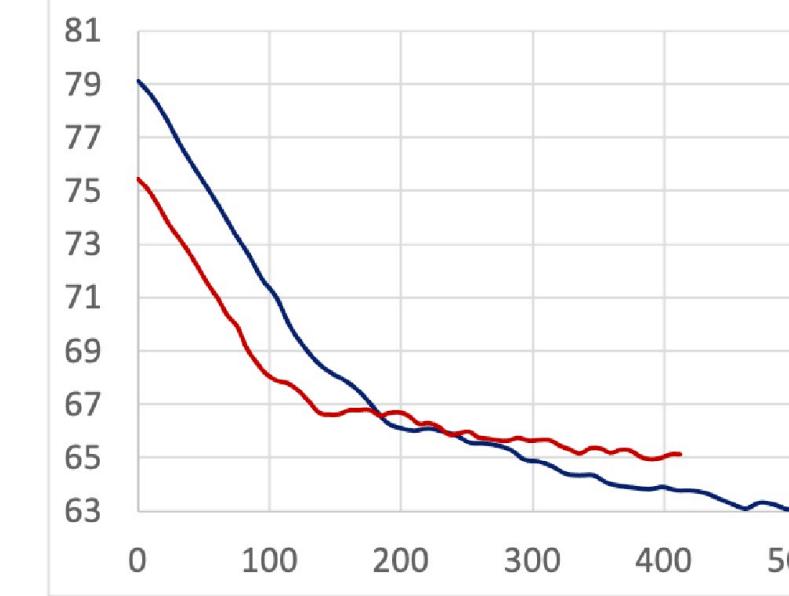
Erroneous Learning: 3/13



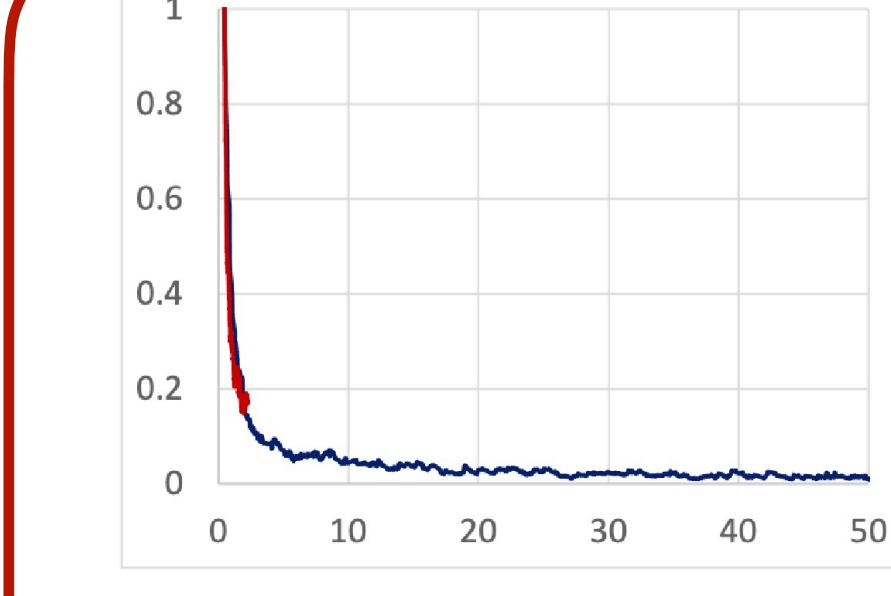
(a) LSTM-MNIST



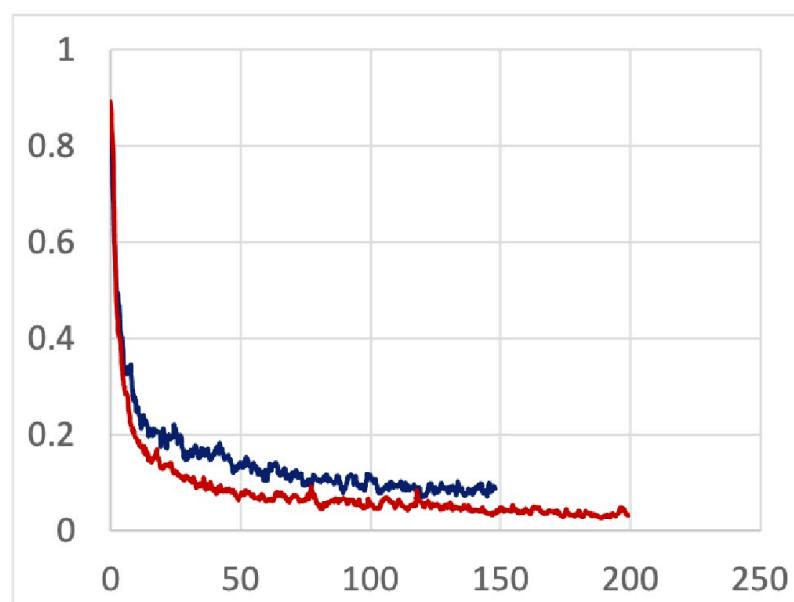
(f) Fashion-MNIST



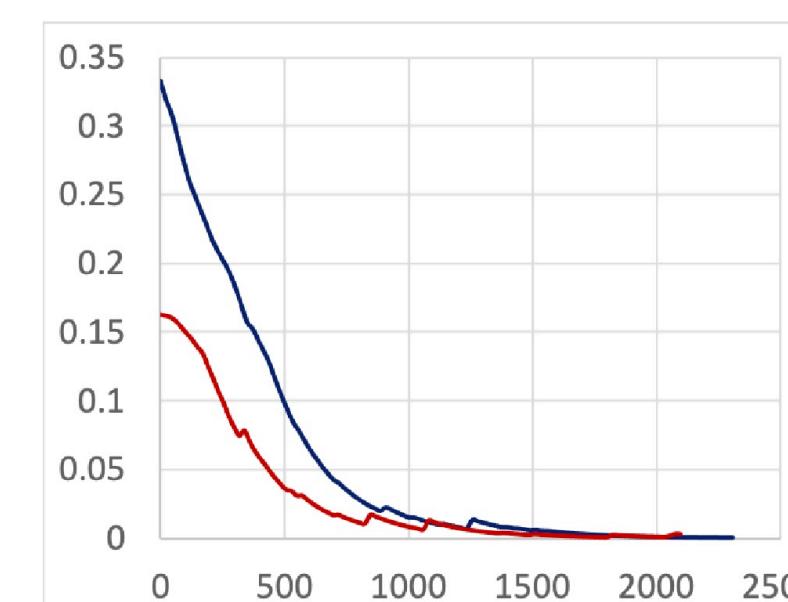
(k) Auto-Encoder



(b) SimpleCNN-GradientTape-2



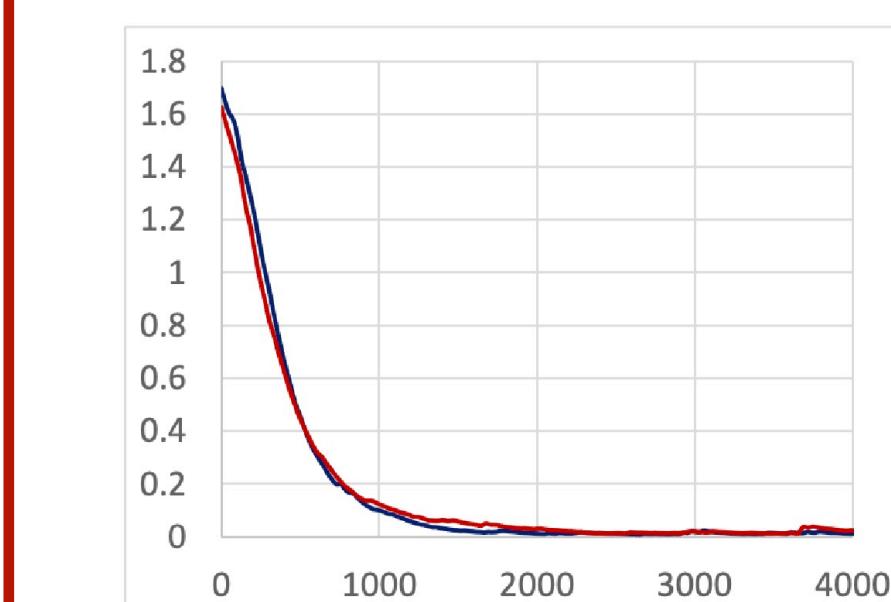
(d) Play-with-MNIST



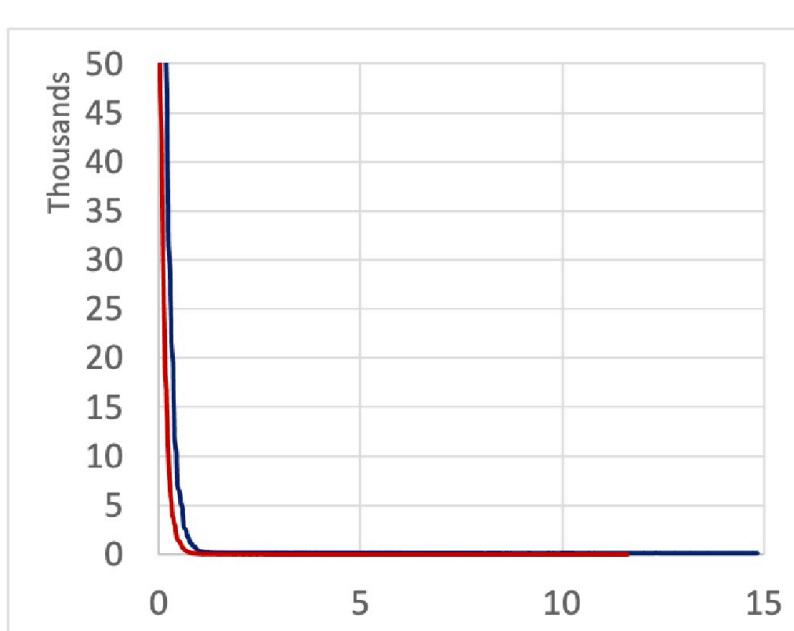
(h) Inception-network



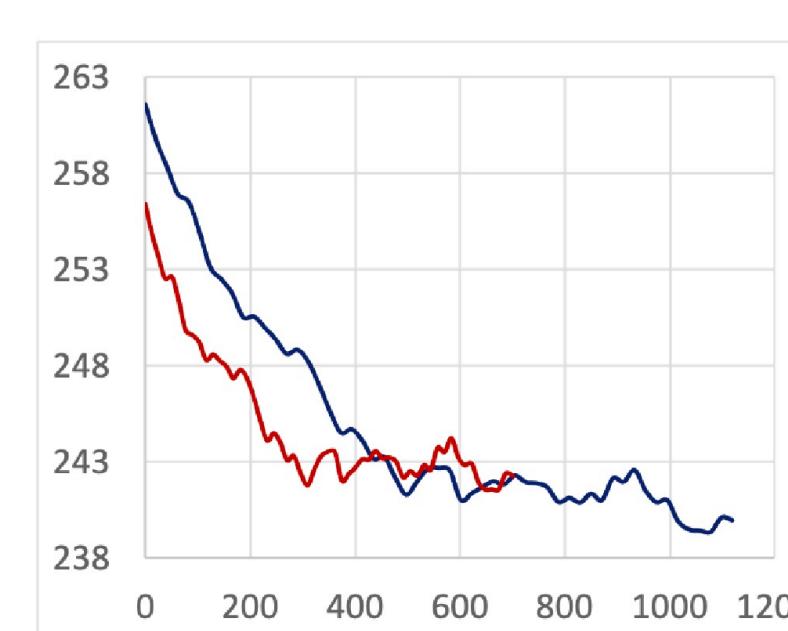
(g) CIFAR10-VGG16



(j) Stacked-LSTM-ColorBot



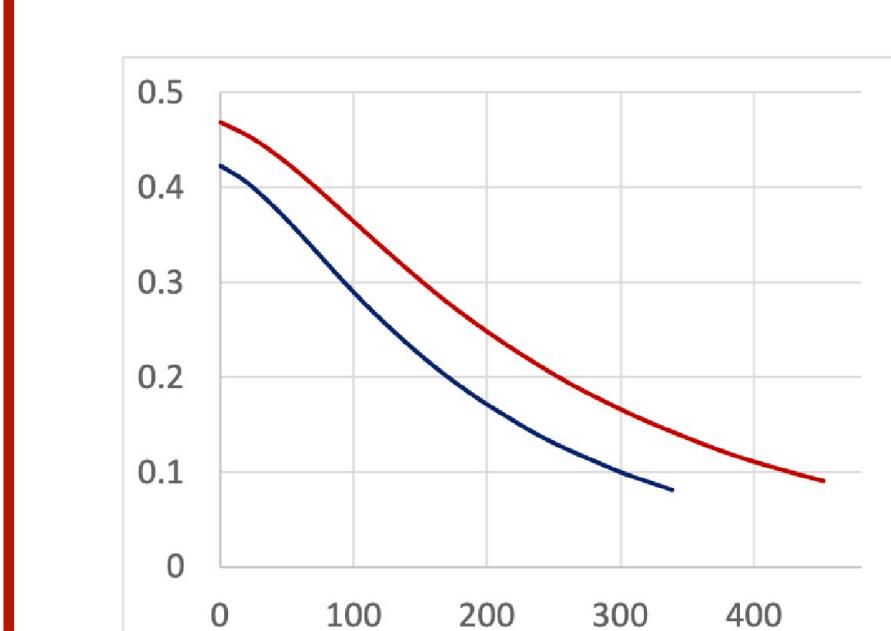
(e) Linear-Regression



(l) Variational-Auto-Encoder

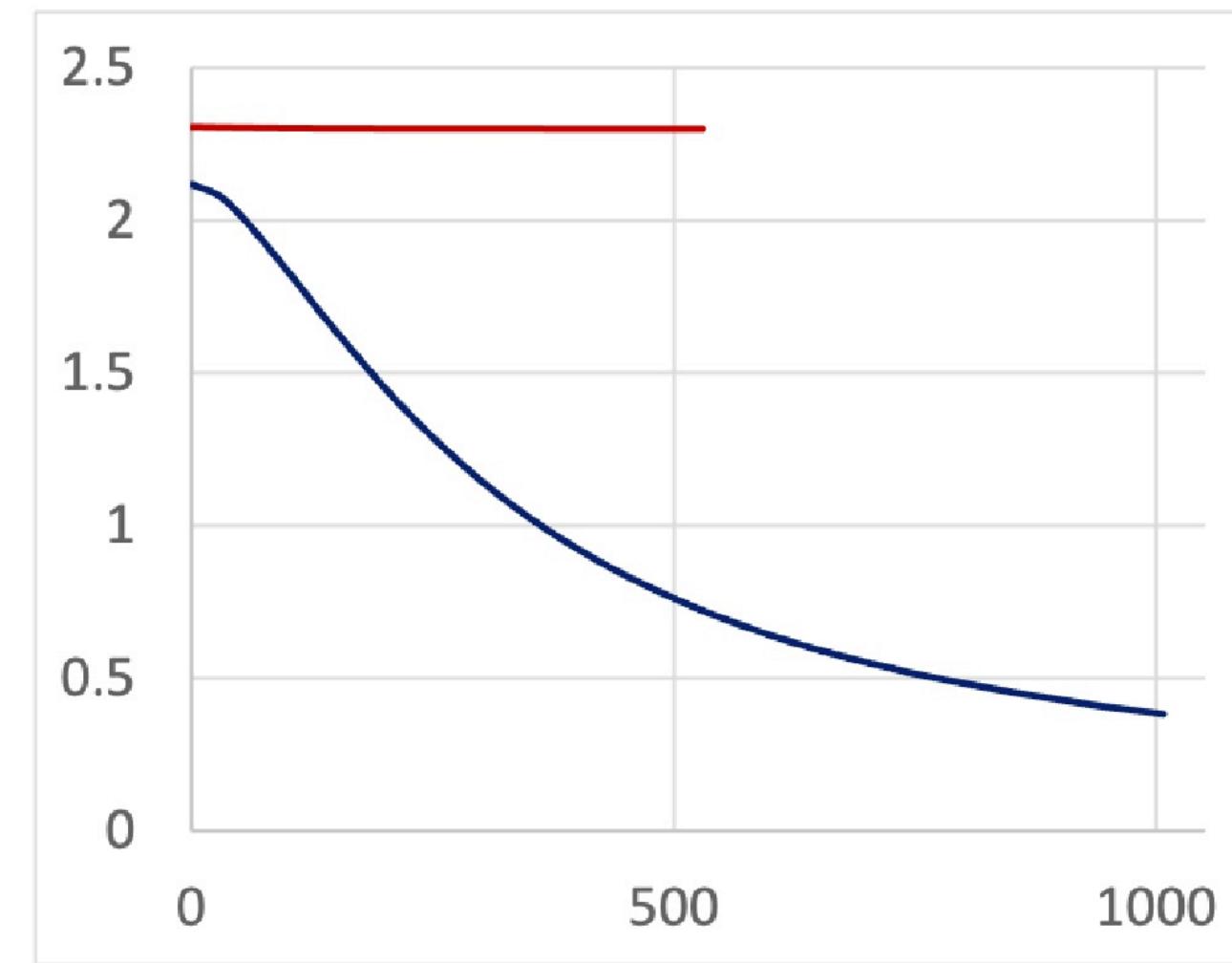


(i) RNN-Sentiment-Analysis

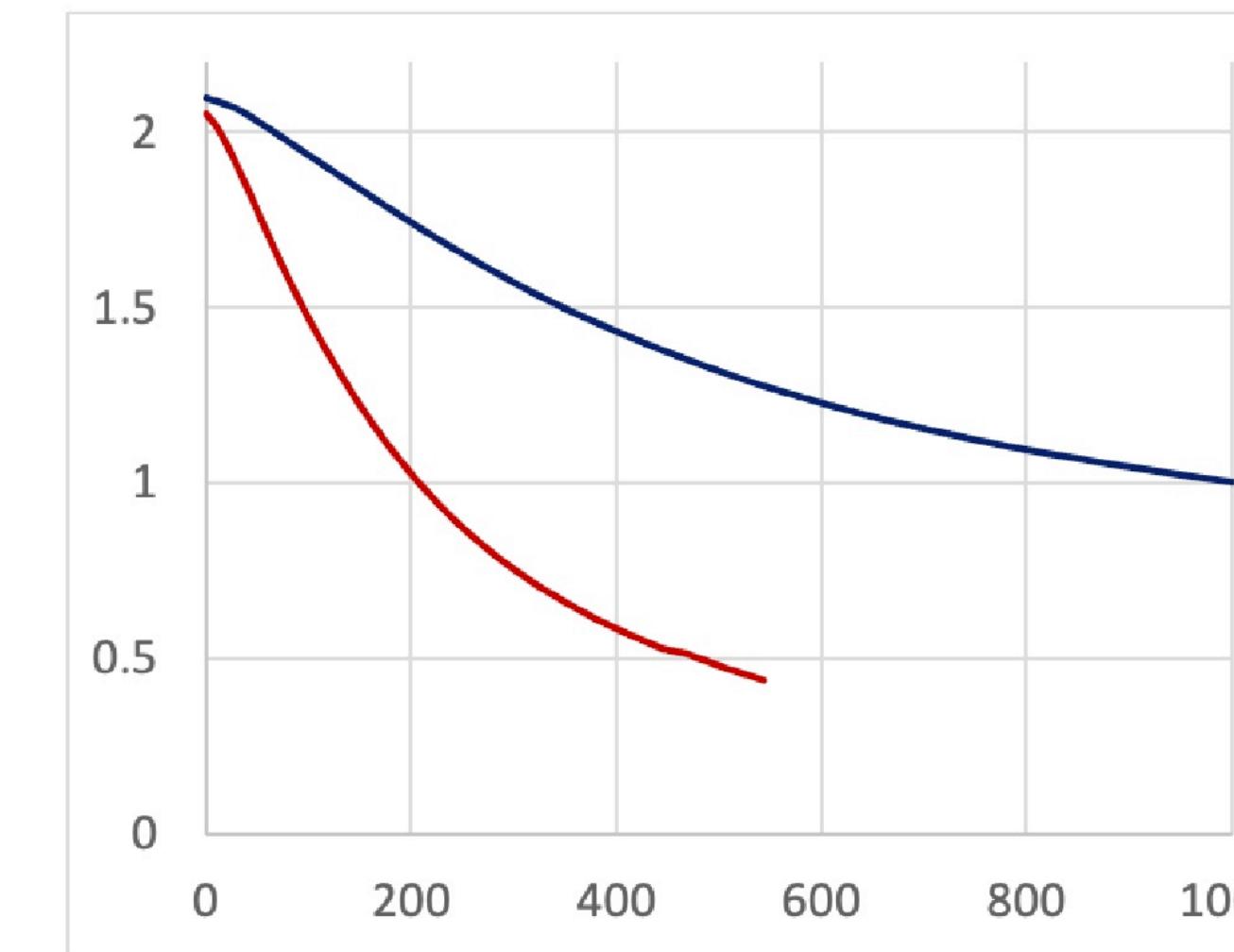


(m) DCGAN

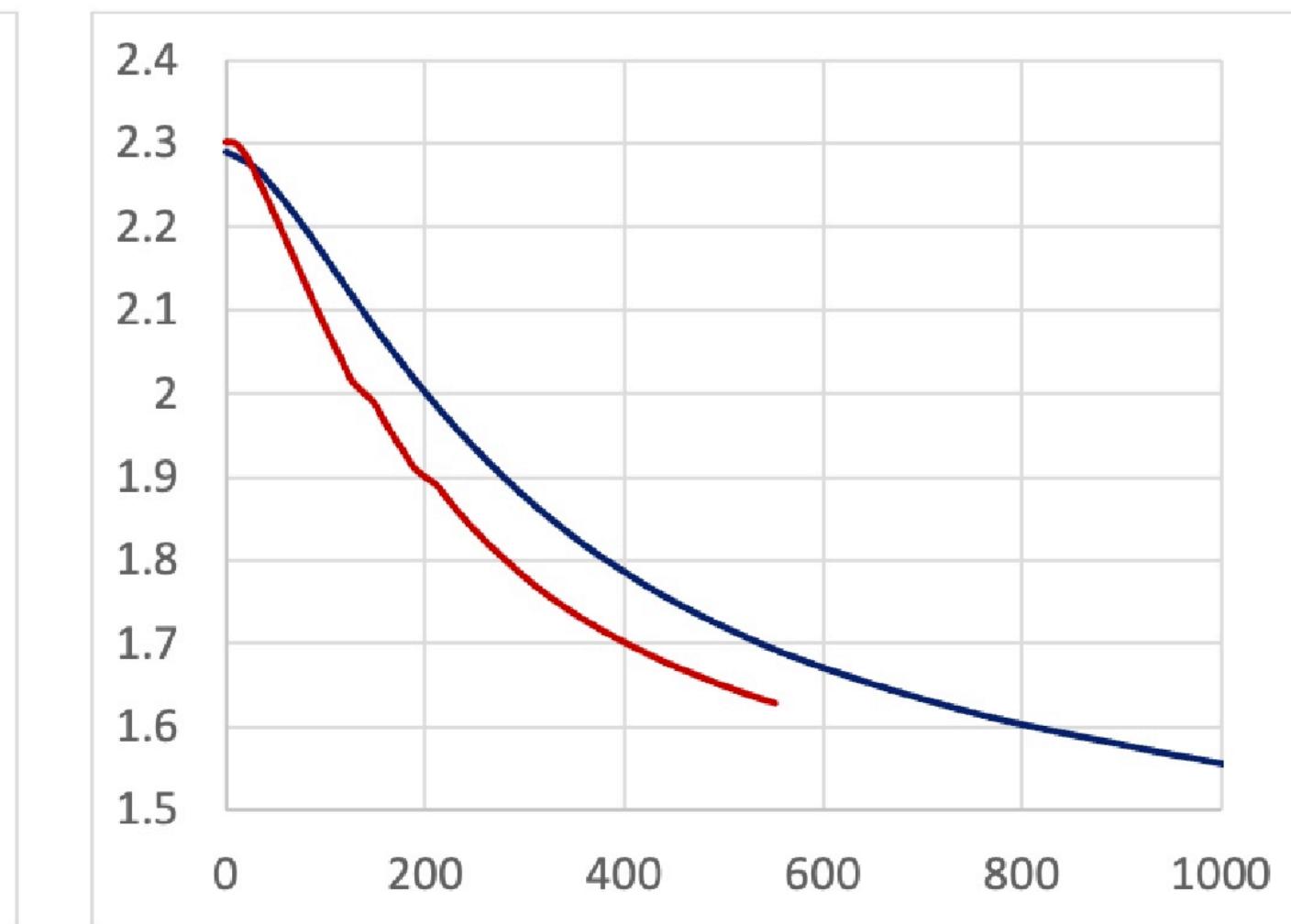
# Evaluation: Analysis on Erroneous Case



(a) VGG-CIFAR10 model with lr=0.001 (default value)



(b) VGG-CIFAR10 model with lr=0.0001



(c) VGG-CIFAR10 model with lr=0.00001

“... GPUs provide many parallel processing units which are ideal for throughput computing. However, the design for graphics pipeline lacks some critical processing capabilities for general-purpose workloads, which may result in lower architecture efficiency on throughput computing workloads.”

- Lee VW et al. Debunking the 100X GPU vs. CPU Myth: An Evaluation of Throughput Computing on CPU and GPU.

# Evaluation: Analysis on Erroneous Case

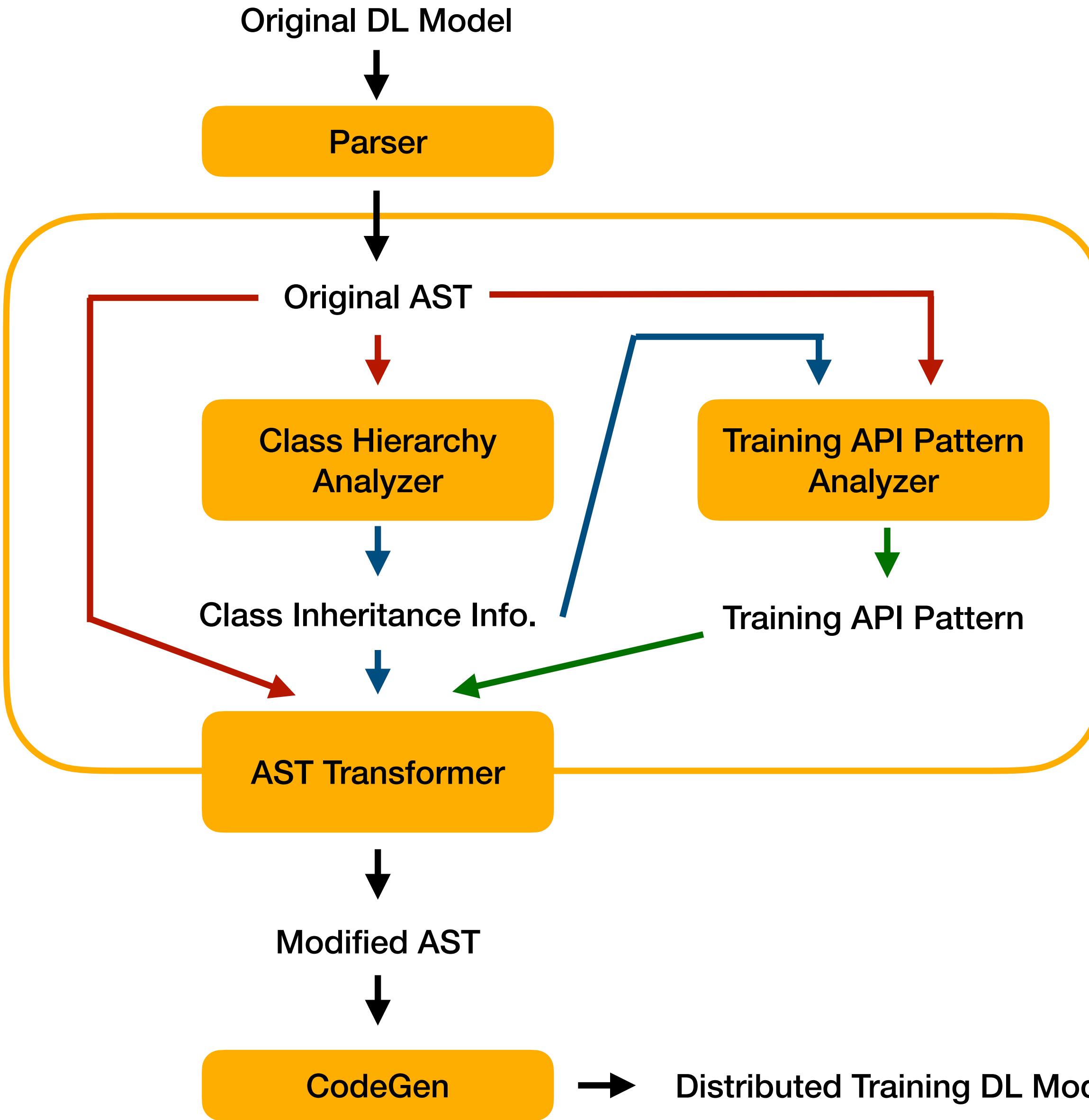
*Developers can automate the rewriting process from single-GPU DL models to distributed DL models, which reduces the time and effort to manually rewrite the model codes.*

*When experimenting with new DL models, developers can measure the training performance in a short time.*

*When using the existing DL model for large-scale products,*

*“... GPUs provide many parallel processing units which are ideal for throughput computing. However, the design for graphics pipeline lacks some critical processing capabilities for general purpose workloads, which may result in lower architecture efficiency on throughput computing workloads.”*

- Lee et al. Debunking the 100X GPU vs. CPU Myth: An Evaluation of Throughput Computing on CPU and GPU.



TF version	Pattern name	Explanation
1.x	Session	Low-level training API using <code>Session</code> instance
1.x	MonitoredSession	Low-level training API using <code>MonitoredSession</code> instance
2.x	GradientTape	Low-level training API using <code>GradientTape</code> instance
2.x	Keras	High-level training API using <code>fit</code> method of <code>keras.models.Model</code> instance

**FIGURE 7** Training API patterns