

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

Monadic Deep Learning

Performing monadic automatic differentiation in parallel

ANONYMOUS AUTHOR(S)

The Java and Scala community has built a very successful big data ecosystem. However, most of neural networks running on it are modeled in dynamically typed programming languages. These dynamically typed deep learning frameworks treat neural networks as differentiable expressions that contain many **trainable variables**, and perform automatic differentiation on those expressions when training them.

Until 2019, none of the learning frameworks in statically typed languages provided the expressive power of traditional frameworks. Their users are not able to use custom algorithms unless creating plenty of boilerplate code for hard-coded back-propagation.

We solved this problem in DeepLearning.scala 2. Our contributions are:

- We discovered a novel approach to perform automatic differentiation in reverse mode for statically typed functions that contain multiple **trainable variables**, and can interoperate freely with the metalanguage.
- We designed a set of monads and monad transformers, which allow users to create monadic expressions that represent dynamic neural networks.
- Along with these monads, we provide some applicative functors, to perform multiple calculations in parallel.

With these features, users of DeepLearning.scala were able to create complex neural networks in an intuitive and concise way, and still maintain type safety.

Additional Key Words and Phrases: type class, path-dependent type, monad, scala

1 BACKGROUND AND INTRODUCTION

Deep neural networks have become the state of the art on many tasks, such as computer vision, game playing, voice recognition and natural language translation.

A neural network is a computation model that transforms the input, into output, by repeated application of tensor operations (including matrix multiplication, tensor resizing, element-wise operations, such as max, +, etc). A “deep” neural network just indicates there are large amount of such transformations.

Additionally, a neural network has additional hidden inputs, called parameters or weights, and deep learning is the task of finding the best weights for a neural network, such that a predefined objective (called loss) is minimized.

Deep learning is mostly done using different variation of gradient descend: we calculate the first order derivative of the loss function with respect to the weight, and update the weight accordingly. The processed is called backpropagation.

Thus, backpropagation[Rumelhart et al. 1985] is the key feature in many deep learning frameworks. Combined with other optimization algorithms[Duchi et al. 2011; Kingma and Ba 2014; Zeiler 2012], deep learning frameworks mutate the values of weights in neural networks during training, producing a model of knowledge learnt from training data.

Backpropagation can be considered as a specialized instance of Automatic Differentiation (AD)[Baydin et al. 2015b]. Many successful Python deep learning frameworks[Google Brain 2017; Neubig et al. 2017; Paszke et al. 2017; Tokui et al. 2015] implement a common set of features of auto differentiation:

50 **Reverse mode** These deep learning frameworks perform reverse mode AD instead of forward
 51 mode, as forward mode AD does not scale well for deep neural networks.

52 **Multiple trainable variable** Neural networks are composed of multiple layers. Each layer con-
 53 tains their own trainable variables. These deep learning frameworks are able to calculate the
 54 derivatives of all trainable variables at once for one training data batch.

55 **Internal DSL**[Fowler 2010] These deep learning frameworks are libraries that provide an Internal
 56 DSL, allowing users to create their differentiable functions in Python or Lua from similar
 57 expressions as creating ordinary non-differentiable functions. Since these frameworks do
 58 not require external language, models created by them can be easily integrated into a larger
 59 application simply with higher-level configurations[Chollet et al. 2015] or ETL (Extract,
 60 Transform and Load) process.

61 However, how to archive those goals in statically typed library is still an open question. Pre-
 62 vious solutions in statically type languages usually requires metaprogramming or compiler-time
 63 translation. In this paper, we present DeepLearning.scala, which is the first statically typed imple-
 64 mentation that achieves all the above goals for deep neural networks. Our approach is based on
 65 some widely available functional programming constructs, thus can be ported to other statically
 66 typed programming languages.

68 2 BASIC CONCEPTS

69 For example, suppose we are building a robot for answering questions in IQ test like this:

70 What is the next number in this sequence:

71 3, 6, 9, ?

72 The answer is 12.

73 In DeepLearning.scala, the robot can be implemented as a guessNextNumber function as follow-
 74 ing¹:

```
76 // Weight initialization
77 val weights: Seq[DoubleWeight] = Stream.continually(DoubleWeight(math.random))
78 val bias: DoubleWeight = DoubleWeight(math.random)
79
80 def guessNextNumber(question: Seq[Double]): DoubleLayer = {
81   // Perform a dot product between question and weights
82   (question zip weights).map {
83     case (element, weight) => element * weight
84   }.reduce(_ + _) + bias
85 }
86 }
```

87 Listing 1. The differentiable matrix multiplication implemented by map/reduce

88
 89
 90 guessNextNumber performs a dot product between question and weights by invoking higher-
 91 order functions map and reduce.

92 Unlike[Chen 2017]’s special tensor type, our tensor can be typed simply as Seq[Double], Seq[
 93 DoubleWeight] or Seq[DoubleLayer].²

94
 95 ¹The code examples from Listing 1 to Listing 7 do not contain necessary import and configurations. For an executable
 96 model backed by ND4J[Skymind 2017b], see [Getting Started documentation on DeepLearning.scala website](#).

97 ²DeepLearning.scala users can use other representations of tensors: (1) For tensors with a statically typed shape, use
 98 [shapeless.Sized](#). For example a 10×20 two-dimensional tensor can be typed as `Sized[Sized[Double, _10], _20]`.

99 The return value of `guessNextNumber`, along with temporary variables in `guessNextNumber`,
 100 are `DoubleLayers`, which are **differentiable expressions**.

101 The weights and bias contains some `DoubleWeight` that are referenced by `guessNextNumber`.
 102 They must be initialized before executing the model. Those are **trainable variables**.

103 From the user’s point of view, both `DoubleLayer` and `DoubleWeight` are opaque types simi-
 104 lar to an ordinary `scala.Double`. Most of the operators of `scala.Double` are also available on
 105 `DoubleLayer` and `DoubleWeight`, except those operators are differentiable. For now, the type sig-
 106 nature of the multiplication operator as can be seen as in Listing 2, and we will reveal the real type
 107 signature of `*` in Section 4.
 108

```

109 trait DoubleLayer {
110     // Scalar multiplication
111     def *(rhs: Double): DoubleLayer
112     def *(rhs: DoubleLayer): DoubleLayer
113     def *(rhs: DoubleWeight): DoubleLayer
114
115     // Element-wise multiplication
116     def *(rhs: INDArray): INDArrayLayer
117     def *(rhs: INDArrayLayer): INDArrayLayer
118     def *(rhs: INDArrayWeight): INDArrayLayer
119 }
120 trait DoubleWeight {
121     // Scalar multiplication
122     def *(rhs: Double): DoubleLayer
123     def *(rhs: DoubleLayer): DoubleLayer
124     def *(rhs: DoubleWeight): DoubleLayer
125
126     // Element-wise multiplication
127     def *(rhs: INDArray): INDArrayLayer
128     def *(rhs: INDArrayLayer): INDArrayLayer
129     def *(rhs: INDArrayWeight): INDArrayLayer
130 }
    
```

131 Listing 2. The hypothetical type signature of multiplication operator for `DoubleLayer` and `DoubleWeight`
 132
 133

134 Table 1 lists `DeepLearning.scala` built-in differentiable types other than `DoubleLayer` and
 135 `DoubleWeight`.
 136

137 Table 1. Built-in Differentiable Types
 138

	non-trainable value	trainable variable	differentiable expression
single-precision scalar	<code>Double</code>	<code>DoubleWeight</code>	<code>DoubleLayer</code>
double-precision scalar	<code>Float</code>	<code>FloatWeight</code>	<code>FloatLayer</code>
vector	<code>INDArray</code>	<code>INDArrayWeight</code>	<code>INDArrayLayer</code>

139
 140
 141 For differentiable tensors, replace vanilla `Double` to `DoubleWeight` or `DoubleLayer`. (2) For GPU-accelerated tensors, use
 142 `INDArray[Skymind 2017b]`. For differentiable tensors, use `INDArrayLayer` or `INDArrayWeight` instead.
 143
 144

148 In addition to differentiable operations, Layers and Weights can be evaluated with a predict
 149 method, thus, the model can predict the next Integer by calling an ordinary function, as shown
 150 in Listing 3. You may notice the blockingAwait suffix appended to predict, because predict
 151 returns a Future[Double]³, which contains the asynchronous task to compute the result. The
 152 actual computation is not performed in predict until blockingAwait is invoked⁴.

```
153
154 val question = Seq(42.0, 43.0, 44.0)
155 println(guessNextNumber(question).predict.blockingAwait)
```

Listing 3. Inference on an untrained model

158 However, guessNextNumber returns an incorrect result because the weights and bias were
 159 randomly initialized, and have not been trained.

160 In order to train them, a loss function is necessary:

```
161
162 def squareLoss(robotAnswer: DoubleLayer, expectedAnswer: Double): DoubleLayer =
163   {
164     val difference: DoubleLayer = robotAnswer - expectedAnswer
165     difference * difference
166   }
```

Listing 4. The differentiable square loss function

169 The above loss function squareLoss determines the squared error between robot’s answer and
 170 the correct answer.

171 Both squareLoss and guessNextNumber are ordinary functions, and can be composed in other
 172 functions:

```
173
174 def linearRegression(question: Seq[Double], expectedAnswer: Double):
175   DoubleLayer = {
176     val robotAnswer = guessNextNumber(question)
177     squareLoss(robotAnswer, expectedAnswer)
178   }
```

Listing 5. A differentiable function to train a linear regression model

181 linearRegression, composed of guessNextNumber and squaredLoss, returns a DoubleLayer
 182 of the loss for a specific question and its expected answer. linearRegression is a linear regression

183 ³ For readers familiar to Haskell, you can understand Future from the corresponding types in Haskell:

- 184 • Future is an opaque type alias of a TryT-transformed UnitContinuation, which is used for asynchronous operations
 185 with the ability of exception handling.
- 186 • “opaque type alias” is similar to the newtype feature in Haskell.
- 187 • TryT provides the ability of exception handling, which is similar to ExceptT monad transformer in Haskell.
- 188 • UnitContinuation is similar to Cont () in Haskell, which is used for asynchronous operations, like an asynchro-
 189 nous version of IO in Haskell.

190 All the above types are general purpose libraries, not parts of DeepLearning.scala. We use those continuation based monadic
 191 data types to archive the ability of parallel and asynchronous execution.

192 ⁴ When predict is used in a real world scenario (e.g. running a neural network in a web service), blockingAwait should
 193 be replaced to flatMap, instead of blocking the current thread, which is an expensive resource. We use this blockingAwait
 194 here because it’s more straightforward for understanding.

197 model with a square loss, and it can be trained as shown Listing 6. The `blockingAwait` is invoked
 198 because `train` returns a `Future[Double]` as well.

```

200 val question1 = Seq(3.0, 4.0, 5.0)
201 val expectedAnswer1 = 6.0
202
203 val question2 = Seq(13.0, 19.0, 25.0)
204 val expectedAnswer2 = 31.0
205
206 for (iteration <- 0 until 500) {
207   linearRegression(question1, expectedAnswer1).train.blockingAwait
208   linearRegression(question2, expectedAnswer2).train.blockingAwait
209 }

```

Listing 6. Training for 500 iterations

213 The weights and bias referenced by `linearRegression` are modified during 500 iterations of
 214 training, by stochastic gradient descent, to minimize the loss returned from `linearRegression`.

215 When weights and bias have been trained to make loss close to zero, `guessNextNumber` should
 216 return values that are very close to the expected answers.

```

218 val question = Seq(42.0, 43.0, 44.0)
219 println(guessNextNumber(question).predict.blockingAwait)
220

```

Listing 7. Inference on a trained model

224 This time, it will print a number close to 45, as the model has finally learned the pattern of
 225 arithmetic progression.

226 The example from Listing 1 to Listing 7 demonstrated some basic concepts in `DeepLearning.scala`.

- 227 • `guessNextNumber`, `squareLoss` and `linearRegression` are **differentiable functions** that
 228 return **differentiable expressions**, which are **computational graph** nodes that can be evaluated
 229 when **training** or **predicting**.
- 230 • **Differentiable expressions** and **trainable variables** can be used as if they are ordinary non-
 231 differentiable values. For example, as shown in Listing 4, you can perform scalar subtraction
 232 and multiplication between `DoubleWeight`, `DoubleLayer` and ordinary `scala.Double`.
- 233 • When **training** a **differentiable expression**, it returns a **Future**, which encapsulates the
 234 side-effect of adjusting **trainable variables** referenced by the **differentiable function**.
- 235 • If a **differentiable function** invokes another **differentiable function**, then **trainable variables**
 236 trained by one **differentiable function** affect another one. For example, when training the
 237 **differentiable function** `linearRegression`, The **trainable variables** weights and bias are
 238 modified, hence `guessNextNumber` automatically gains the ability to **predict** correct an-
 239 swers.

241 3 DYNAMIC NEURAL NETWORKS

242 `DeepLearning.scala` supports dynamic neural networks. It means that the control flow of a neu-
 243 ral network can differ according to values of internal nodes of the computational graph, when
 244 processing a specific batch of input. This is the key feature of recent deep learning frameworks

246 like PyTorch[Paszke et al. 2017] or Chainer[Tokui et al. 2015]. Especially, dynamic deep neural
 247 networks can be more efficient by skipping part of the model[Liu and Deng 2017].

248 In this section, we will present how to create a simple dynamic neural network in DeepLearn-
 249 ing.scala, which can be considered as a simplified version of outrageously large neural net-
 250 works[Shazeer et al. 2017].

251 Suppose we have two sub-neural networks, leftSubnet and rightSubnet (Listing 8)⁵. We want
 252 to build a “gated” network, which conditionally runs either leftSubnet or rightSubnet for a
 253 special input.

```
255 def leftSubnet(input: INDArraryLayer): INDArraryLayer
256 def rightSubnet(input: INDArraryLayer): INDArraryLayer
```

257 Listing 8. Predefined sub-networks

258
 259
 260 Which sub-network is selected for the input should be determined by the gate network, which
 261 returns a pair of differentiable double expressions that indicate the preferences between leftSubnet
 262 and rightSubnet, shown in (Listing 9).

```
264 def gate(input: INDArraryLayer): (DoubleLayer, DoubleLayer)
```

265 Listing 9. Predefined gate network

266
 267
 268 The differentiable operations on DoubleLayers and INDArraryLayers form a differentiable em-
 269 bedded DSL inside the meta-language Scala. Thus, the concept of “gated” neural network can be
 270 considered as a conditional control flow in the differentiable DSL, and the values of nodes of a
 271 gated neural network are simply some let bindings in the DSL.

272 The control flow of the gated network that we want to build is described in Function GatedNet-
 273 work.

274 **Function** GatedNetwork

275 **Input:** Features extracted by preceding layers
 276 **Output:** Features passed to succeeding layers
 277 scores \leftarrow gate(*Input*);
 278 **if** *score of left sub-network* > *score of right sub-network* **then**
 279 | **return** *score of left sub-network* \times leftSubnet(*Input*);
 280 **else**
 281 | **return** *score of right sub-network* \times rightSubnet(*Input*);
 282 **end**

283
 284 In DeepLearning.scala, there are three different approaches to implement the gated network.
 285 Examples of these approaches are introduced in following Section 3.1, Section 3.2, and Section 3.3.

286 3.1 Eager Execution (bad)

287
 288 An obvious approach to create the gated network is to eagerly execute the gate, shown in Listing 10:

289 There are three sub-networks in the naiveGatedNet function. The gate returns a pair of
 290 DoubleLayers. By blocking await the prediction result, we get two Doubles, which can be used

291 ⁵ For performance purpose, instead of Seq[DoubleLayer], we use INDArraryLayer as the type of **differentiable expression**
 292 backed by ND4J. INDArraryLayer supports differentiable version of most operations that ND4J’s n-dimensional array
 293 INDArrary supports.

```

295 def naiveGatedNet(input: INDArrayLayer): INDArrayLayer = {
296   val scores = gate(input)
297   if (scores._1.predict.blockingAwait > scores._2.predict.blockingAwait) {
298     scores._1 * leftSubnet(input)
299   } else {
300     scores._2 * rightSubnet(input)
301   }
302 }

```

Listing 10. The eager execution implementation of gated network

to determine which sub-network is preferred between `leftSubnet` and `rightSubnet`. The chosen sub-network will be multiplied with the value returned by the gate in order to enabling backpropagation on the gate.

However, there is a performance issue in the `naiveGatedNet`.

In `DeepLearning.scala`, all differentiable expressions, including the scalar `DoubleLayer` and vectorized `INDArrayLayer`, contain some lazily evaluated differentiable **computational graph** nodes, which will not be executed until a `predict` or `train` task is performed in a `blockingAwait` or `onComplete` call.

So, the two `predict.blockingAwait` calls in the `if` will execute the **computational graph** in gate twice. Also the **computational graph** in `naiveGatedNet` will be executed when users call `predict` or `train` in the future. Even worse, `input` contains a **computational graph**, too. Along with `gate`, it will be evaluated three times, which may contain a complex future extracting process.

3.2 Monadic Control Flow (good)

Ideally, the calls to `predict` should be avoided in differentiable functions. The recommended approach to create a dynamic neural network is using `forward`, which returns a monadic value of `Do[Tape[Data, Delta]]`, which can be used in a monadic control flow via Scalaz[Yoshida 2017]’s type classes[Oliveira et al. 2010] `Monad` and `Applicative`.

Listing 11 shows the monadic control flow of a gated network. The gated network is built from the monadic expression `gatedForward`, which contains some `forward` calls, which are asynchronous operations (or `Do`) that produce Wengert list records (or `Tapes`). The implementation details of `Do` and `Tape` will be discussed in Section 5. For now, we only need to know that `Do` is a monadic data type that supports `flatMap`.

By `flatMap`ing those `forward` operations together, we built the entire monadic control flow `gatedForward` for the gated network.

The `monadicGatedNet` represents a dynamic neural network, since each `forward` operation is started after its previous `forward` is done. This behavior allows for dynamically determining succeeding operations according to results of previous `forward` operations, as shown in the `if` clause in Listing 11.

However, `flatMap` prevents additional optimization, too. `scores._2.forward` have to wait for `scores._1.forward`’s result, even if the two operations are logically independent.

3.3 Parallel Applicative Control Flow + Sequential Monadic Control Flow (better)

Ideally, the independent operations `scores._1.forward` and `scores._2.forward` should run in parallel. This can be done by tagging `Do` as `Parallel`, and use `scalaz.Applicative.tuple2` [McBride and Paterson 2008] instead of `flatMap` (Listing 12).

```

344 def monadicGatedNet(input: INDArrayLayer): INDArrayLayer = {
345   val scores = gate(input)
346   val gatedForward: Do[Tape[INDArray, INDArray]] = {
347     scores._1.forward.flatMap { tape1: Tape[Double, Double] =>
348       scores._2.forward.flatMap { tape2: Tape[Double, Double] =>
349         if (tape1.data > tape2.data) {
350           (scores._1 * leftSubnet(input)).forward
351         } else {
352           (scores._2 * rightSubnet(input)).forward
353         }
354       }
355     }
356   }
357   INDArrayLayer(gatedForward)
358 }

```

Listing 11. Monadic gated network

```

362 def applicativeMonadicGatedNet(input: INDArrayLayer): INDArrayLayer = {
363   val scores = gate(input)
364   val parallelForward1: ParallelDo[Tape[Double, Double]] = {
365     Parallel(scores._1.forward)
366   }
367   val parallelForward2: ParallelDo[Tape[Double, Double]] = {
368     Parallel(scores._2.forward)
369   }
370   val Parallel(stage1) = {
371     parallelForward1.tuple2(parallelForward2)
372   }
373   def stage2(tapes: (Tape[Double, Double], Tape[Double, Double])) = {
374     if (tapes._1.data > tapes._2.data) {
375       (scores._1 * leftSubnet(input)).forward
376     } else {
377       (scores._2 * rightSubnet(input)).forward
378     }
379   }
380
381   val gatedForward = stage1.flatMap(stage2)
382   INDArrayLayer(gatedForward)
383 }

```

Listing 12. Applicative + monadic gated network

388 This `applicativeMonadicGatedNet` takes both advantages from applicative functors and mon-
389 ads. The entire control flow is a `flatMap` sequentially composed of two stages. In `stage1`, there is
390 a `tuple2` composed of `scores._1.forward` and `scores._2.forward` in parallel. Then, in `stage2`,
391 the succeeding operation is dynamically determined according to `tapes`, the result of `stage1`.

392

The parallel applicative operation is also the default behavior for all built-in vector binary operators. Listing 13 shows some simple expressions that will be executed in parallel.

```

396 def parallelByDefault(a: INDArrayLayer, b: INDArrayLayer, c: INDArrayLayer, d:
397     INDArrayLayer): INDArrayLayer = {
398     a * b + c * d
399 }

```

Listing 13. By default, $a * b$ and $c * d$ will be executed in parallel because they are independent

By combining both applicative functors and monads, `DeepLearning.scala` supports dynamic neural network and still allows the independent parts of the neural network to run in parallel. In addition, the backward pass of differentiable functions built from parallel applicative functors or built-in vector binary operators will be executed in parallel, too.

3.4 Direct style DSL for Applicative and Monadic Control Flow (best)

In the previous section, we had presented a dynamic neural network executed in parallel. However, the usage of `flatMap` and `Parallel`-tagged types may scare algorithm authors who are not familiar with monadic programming. Ideally, the code written by those people should look straightforward and has the same structure in pseudo-code `GatedNetwork` or Listing 10, and still gain the benefits of Listing 12.

The goal can be achieved by transforming the direct style code into monadic and applicative code at compile-time. We created an DSL with the help of the `!`-notation provided by `Dsl.scala` [Yang 2017], which provides Scala compiler plugins to perform the necessary compiler-time transformation.

As shown in Listing 14, the `!`-notation “extracts” the `Double` values from a pair of `DoubleLayers` in parallel. Those `Doubles` are ordinary non-differentiable Scala types that will not backpropagate, and can be used in ordinary Scala control flow expression like `if`.

```

421 def dslGatedNet(input: INDArrayLayer): INDArrayLayer = {
422     val scoreLayers: (DoubleLayer, DoubleLayer) = gate(input)
423     val scores: (Double, Double) = !scoreLayers
424     if (scores._1 > scores._2) {
425         scoreLayers._1 * leftSubnet(input)
426     } else {
427         scoreLayers._2 * rightSubnet(input)
428     }
429 }

```

Listing 14. `Dsl.scala` powered direct style gated network

Generally, a `!`-notation on a `Layer` will generate a monadic `flatMap` call, to extract the value of forward pass of the `Layer`; a `!`-notation on a tuple of `Layers` will generate some applicative `<*>` and `map` calls, to extract a tuple of values of forward pass of those `Layers`, in parallel. Thus, the actually code generated by `Dsl.scala`’s compiler plugins for Listing 14 is similar to Listing 12.

4 AD HOC POLYMORPHIC DIFFERENTIABLE FUNCTIONS

In section 2, we had presented the hypothetical differentiable types of multiplication for `DoubleLayer` and `DoubleWeight`. However, the method overloading approach shown in Listing 2 is too verbose,

442 and requires a lot of boilerplate code. In this section, we will present an approach to create custom
 443 **differentiable functions**, without such methods overloading.

444 A **differentiable function** is a neural network. Ideally, a differentiable function should be an ad
 445 hoc polymorphic function that accepts heterogeneous types of parameters, including:

- 446 • A vanilla vector input. i.e. `INDArray`.
- 447 • **Differentiable expressions** of hidden states produced by any previous neural network layers,
 448 i.e. any `INDArrayLayers` regardless of the prefixes.
- 449 • **Trainable variables** in the case of activation maximization technique[Erhan et al. 2009].i.e.
 450 any `INDArrayWeights` regardless of the prefixes.
- 451 • Other user defined differentiable type.

452 Table 1 shows nine types that have built-in `DeepLearning` type classes.

453 This can be achieved with the help of[Gurnell 2017]’s Scala encoding of dependent-type type
 454 class. We defined a `DeepLearning` (Listing 15⁶) type class that witnesses any supported expressions
 455 including **differentiable expressions**, **trainable variables**, or vanilla non-differentiable types. The
 456 users can create type aliases to restrict the types of state during forward pass and backward pass as
 457 shown in Listing 16.
 458

```

459 trait DeepLearning[Differentiable] {
460   /** The type of result calculated during forward pass */
461   type Data
462
463   /** The type of derivative during backward pass */
464   type Delta
465
466   def forward(differentiable: Differentiable): Do[Tape[Data, Delta]]
467   // Other auxiliary methods is omitted
468 }
  
```

469
 470 Listing 15. The dependent-type type `DeepLearning`

```

473 type INDArrayExpression[Expression] = DeepLearning[Expression] {
474   /** The type of result calculated during forward pass */
475   type Data = INDArray
476
477   /** The type of derivative during backward pass */
478   type Delta = INDArray
479 }
  
```

480
 481 Listing 16. A type class alias that witnesses dense vector expressions

482
 483 By using `INDArrayExpression` as a context bound, we can create a polymorphic differentiable
 484 function that accepts any vector expression.

485 ⁶ For readers who are more familiar with Idris, there is a corresponding notation in Idris:

```

486 interface DeepLearning Differentiable where
487   Data : Type
488   Delta : Type
489   forward : Differentiable -> Do (Tape Data Delta)
  
```

```

491 def polymorphicDifferentiableFunction[A: INDArrayExpression, B:
492     INDArrayExpression, C: INDArrayExpression, D: INDArrayExpression](a: A, b:
493     B, c: C, d: D): INDArrayLayer = {
494     a * b + c * d
495 }

```

Listing 17. A polymorphic differentiable function

Listing 17 is similar to Listing 13, except each argument of `polymorphicDifferentiableFunction` accepts `INDArray`, `INDArrayWeight` or `INDArrayLayer` respectively, not only `INDArrayLayer`.

Built-in operations including arithmetic operations, `max`, and `dot` are polymorphic differentiable functions defined in this approach, too, which can be used in user-defined polymorphic **differentiable functions**.

5 IMPLEMENTATION

In this section, we will introduce the internal data structure used in `DeepLearning.scala` to perform AD.

- For ease of understanding, Section 5.1 starts from a simple dual number implementation `DualNumber`, which was known as an approach to perform forward mode AD for scalar values.
- Section 5.2 introduces our variation of dual number `ClosureBasedDualNumber`, which supports tree-structured reverse mode AD (aka backpropagation) for multiple **trainable variables**.
- Section 5.3 shows the actual data type `Tape` in `DeepLearning.scala`, which is generalized to not only scalar types, but also vector types and any other differentiable types.
- Section 5.4 discovers the monadic control flow `Do`, which manages the life cycle of `Tapes`, sharing `Tapes` for common **computational graph** nodes, allowing for an arbitrary DAG (Directed Acyclic Graph)-structured **computational graph**.
- Section 5.5 summarizes the entire execution process during a training iteration, showing how the user-defined **differentiable functions** get executed through internal mechanisms `Do` and `Tape`.

5.1 Ordinary Dual Number

Our approach for reverse mode AD uses a data structure similar to traditional forward mode AD, with only a few changes.

Forward mode AD can be considered as computation on dual number. For example, dual number for scalar types can be implemented as Listing 18:

```

528 type Data = Double
529 type PartialDelta = Double
530 case class DualNumber(data: Data, delta: PartialDelta)

```

Listing 18. Dual number for forward mode AD

Arithmetic operations on those dual number can be implemented as Listing 19:

5.2 Monadic Closure-based Dual Number

However, it is hard to type this approach if we want to support multiple **trainable variables**, with the number unknown before runtime. `PartialDelta` in Listing 18 represents the partial derivative

```

540 object DualNumber {
541     def plus(left: DualNumber, right: DualNumber): DualNumber = {
542         DualNumber(left.data + right.data, left.delta + right.delta)
543     }
544     def multiply(left: DualNumber, right: DualNumber): DualNumber = {
545         DualNumber(left.data * right.data, left.data * right.delta + right.data *
546             left.delta)
547     }
548 }

```

Listing 19. Arithmetic operations on dual number

550

of **trainable variables**. In AD tools that support only one **trainable variable**, the **trainable variable** is usually forced to be the input. Hence `PartialDelta` is the input type for those AD tools. This assumption is broken for our case, since our `delta` type of a specific `DualNumber` must contain derivatives for all **trainable variables** that were used to produce the `DualNumber`, not only the partial derivative of input. As a result, the type of `delta` varies when the number of **trainable variables** grows.

To type-check the `delta`, considering the only usage of the `delta` in a neural network is in updating **trainable variables** in a gradient descent based optimization algorithm, we can replace `PartialDelta` to an `UpdateWeights` closure as shown in Listing 20.

```

562 type Data = Double
563 case class ClosureBasedDualNumber(data: Data, backward: UpdateWeights)

```

Listing 20. Replacing `PartialDelta` to a closure

566

Unlike `PartialDelta`, `UpdateWeights` is not a number that supports native arithmetic operations, thus we have to replace native arithmetic operations on `PartialDelta` to some static functions on `UpdateWeights` when implement arithmetic operations for the new dual number, as shown in Listing 21

The only question remaining is to implement the `UpdateWeights`, to make its behavior be equivalent to the original dual number implementation.

Mathematically, the `UpdateWeights` type in a dual number should form any vector space, i.e. the `UpdateWeights` closure itself must support addition and scalar multiplication operations.

Our approach is making `UpdateWeights` be a function type that contains side-effects to update **trainable variables**. Thus the addition operation for closures can be defined as (1).

$$(f_0 + f_1)(x) = f_0(x) + f_1(x) \quad (1)$$

And the scalar multiplication operation for closures is defined as (2):

$$(x_0 f)(x_1) = f(x_0 x_1) \quad (2)$$

The above definition of arithmetic operations can be implemented in monadic data types as shown in Listing 22.

`UpdateWeights`, as a replacement to original `PartialDelta`, is a closure able to update derivatives for all weights with a coefficient (the `Double` parameter). `|+|` is the append operation of `Semigroup`, which could be any cumulative data type.

588

```

589
590 object ClosureBasedDualNumber {
591   def plus(left: ClosureBasedDualNumber, right: ClosureBasedDualNumber):
592     ClosureBasedDualNumber = {
593       ClosureBasedDualNumber(left.data + right.data, UpdateWeights.plus(left.
594         backward(), right.backward()))
595     }
596   def multiply(left: ClosureBasedDualNumber, right: ClosureBasedDualNumber):
597     ClosureBasedDualNumber = {
598       ClosureBasedDualNumber(
599         left.data * right.data,
600         UpdateWeights.multiply(left.data, right.backward()) + UpdateWeights.
601         multiply(right.data, left.backward()))
602     }
603 }

```

Listing 21. Replacing operations on PartialDelta to custom functions for UpdateWeights

```

607 type UpdateWeights = Do[Double] => SideEffects
608 object UpdateWeights {
609   /**  $(f_0 + f_1)(x) = f_0(x) + f_1(x)$  */
610   def plus(f0: UpdateWeights, f1: UpdateWeights) = { doX: Do[Double] =>
611     f0(doX) |+| f1(doX)
612   }
613
614   /**  $(x_0 f)(x_1) = f(x_0 x_1)$  */
615   def multiply(x0: Double, f: UpdateWeights) = { doX1: Do[Double] =>
616     f(doX1.map(x0 * _))
617   }
618 }

```

Listing 22. Arithmetic operations for the closure that contains side-effects

Also note that the parameter is a monadic data type `Do` that encapsulates the computation of derivative. Unlike strictly evaluated values, `Do` is an operation only evaluated when needed.

In `DeepLearning.scala`, our `SideEffects` is based on the asynchronous operation type `UnitContinuation`. Our built-in differentiable operations execute the independent parts of backward pass in parallel with the help of `Applicative` type class instances of `UnitContinuation`.

```

628
629 type SideEffects = UnitContinuation[Unit]

```

Listing 23. Monadic side-effects

`UnitContinuation[A]` is an opaque alias [Osheim and Cantero 2017] of `(A => Trampoline[Unit]) => Trampoline[Unit]`, implemented in a separate library at [future.scala](#). It is used in `DeepLearning.scala` as a monadic data type for encapsulating side effects in stack-safe asynchronous programming.

638 The SideEffects for neural networks conform to the associative law because the only side effects
 639 is updating **trainable variables**. Thus, our UpdateWeights.plus and UpdateWeights.multiply
 640 are equivalent to the operations on strictly evaluated scalar value PartialDelta in forward mode
 641 AD.

642 Since UpdateWeights is a closure with side effects, a **trainable variable** can be represented as a
 643 tuple of a mutable value and the action to modify the mutable value.

```
645 def createTrainableVariable(initialValue: Double, learningRate: Double):
646     ClosureBasedDualNumber = {
647     var data = initialValue
648     val backward: UpdateWeights = { doDelta: Do[Double] =>
649         val sideEffects: Do[Unit] = doDelta.map { delta =>
650             data -= learningRate * delta
651         }
652         convertDoToContinuation(sideEffects)
653     }
654     ClosureBasedDualNumber(data, backward)
655 }
```

Listing 24. Create a dual number for a **trainable variable**

659 In Listing 24, the **trainable variable** is trained by a fixed learning rate to simplify the hyperpa-
 660 rameters of optimization algorithms. The sideEffects is the only code in the entire code base
 661 that performs native side effects. The actual DeepLearning.scala implementation has a more so-
 662 phisticated mechanism to customize the implementation of side effects, allowing other optimizers
 663 and custom hyperparameters.

664 Similar to **trainable variables**, a non-trainable value can be represented as a tuple of the value
 665 and a no-op closure as shown in Listing 25.

```
667 def createLiteral(data: Double): ClosureBasedDualNumber = {
668     val backward = { doDelta: Do[Double] =>
669         UnitContinuation.now(())
670     }
671     ClosureBasedDualNumber(data, backward)
672 }
```

Listing 25. Create a dual number for a non-trainable value

676 Because delta is an action instead of pre-evaluated value, the implementation of backward for
 677 a non-trainable value can entirely avoid executing any unnecessary computation in doDelta.

678 Finally, we can create a differentiable function as shown in Listing 26, whose leaf nodes are
 679 createTrainableVariable and createLiteral, and internal nodes are arithmetic operations in
 680 Listing 21.

681 The **computational graph** of computationalTree is shown in Figure 1. Note that the arrow
 682 direction denotes the dependency between expressions, from arrow tail to arrow head, which is
 683 the reverse of the direction of data flow.

684 The closure-based dual number y1 has a closure backward, which returns SideEffects that
 685 recursively change all **trainable variables** referenced by the closure.

686

```

687 val w0 = createTrainableVariable(math.random, 0.001)
688 val w1 = createTrainableVariable(math.random, 0.001)
689
690 def computationalTree(x: ClosureBasedDualNumber) = {
691   val y0 = ClosureBasedDualNumber.multiply(x, w0)
692   val y1 = ClosureBasedDualNumber.multiply(y0, w1)
693   y1
694 }

```

Listing 26. A tree-structured differentiable function

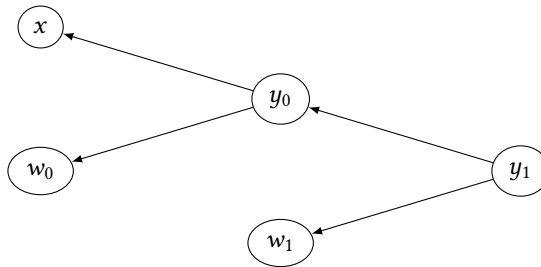


Fig. 1. A tree-structured computational graph

Note that backward itself does not perform any side effects. It just collects all side effects into a UnitContinuation[Unit]. Figure 2 shows how the side effects of updating trainable variables are collected.

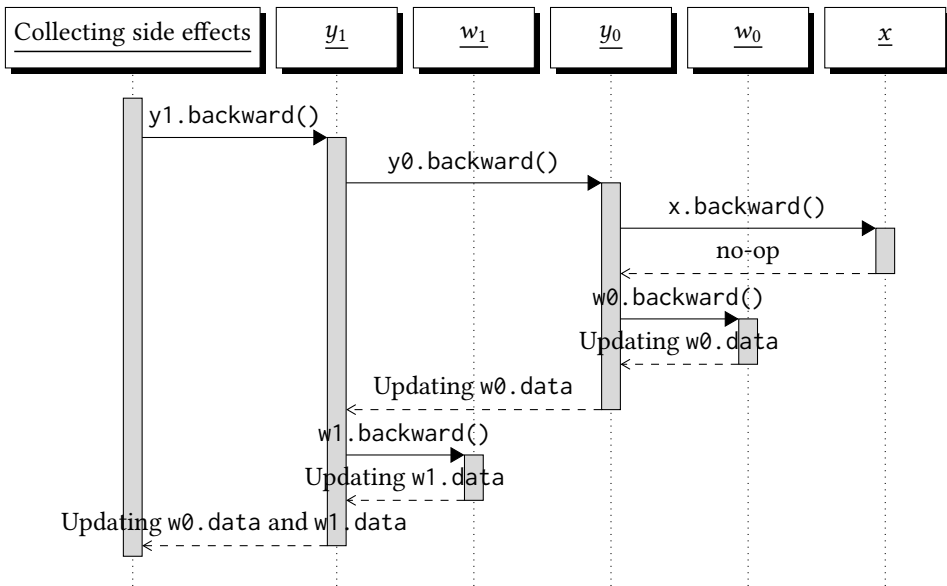


Fig. 2. Backpropagation for a tree-structured computational graph

Finally, the collected side effects of `UnitContinuation[Unit]` returned from `y1.backward` can be performed by a `blockingAwait` or `onComplete` call.

5.3 Generic Tape

This closed-based monadic dual number can be generalized to any linear spaces, not only scalar types such as `Double`, but also n-dimensional arrays.

The dual number type that we actually defined in `DeepLearning.scala` is `Tape`, a generic version of `ClosureBasedDualNumber` in Listing 20. We replaced `ClosureBasedDualNumber`'s hard-coded `Double` to type parameters `Data` and `Delta`, as shown in Listing 27.

```
final case class Tape[+Data, -Delta](
  data: Data,
  backward: Do[Delta] => UnitContinuation[Unit]
)
```

Listing 27. Generic closed-based monadic dual number

`Data` and `Delta` are usually the same, but they can also be different types. For example, you can create a type whose `Data` is a dense n-dimensional array and whose `Delta` is a pair of index and scalar, representing a dense tensor that sparsely updates.

This data structure is similar to a Wengert list in traditional reverse mode AD, except our tape is a tree of closures instead of a list.

5.4 Reference Counted Tape

Although the closed-based dual number approach from Listing 20 to Listing 27 supports multiple trainable variables, the closure-based computation has a performance issue in the case of diamond dependencies.

Listing 28 shows a `differentiable function` `diamondDependentComputationalGraph` that contains diamond dependencies to some `differentiable expressions` or `trainable variables`.

```
val w = createTrainableVariable(math.random, 0.001)
def diamondDependentComputationalGraph(x: ClosureBasedDualNumber) = {
  val y0 = ClosureBasedDualNumber.multiply(x, w)
  val y1 = ClosureBasedDualNumber.multiply(y0, y0)
  val y2 = ClosureBasedDualNumber.multiply(y1, y1)
  y2
}
```

Listing 28. A diamond dependent `differentiable function`

The `computational graph` of `diamondDependentComputationalGraph` are shown in Figure 3.

When `y2.backward` is invoked, in order to collect side effects of `y2`'s dependencies, `y1.backward` will be invoked, twice, and each `y1.backward` call will triggers two `y0.backward` calls. As a result, for each iteration of backpropagation, `y0.backward`, `w.backward` and `x.backward` are invoked four times, respectively.

The process in `y2.backward` is shown in Figure 4.

Generally, given n levels of nested diamond dependencies, the computational complexity is $O(2^n)$, which is unacceptable for neural networks that may share common `differentiable expressions`.

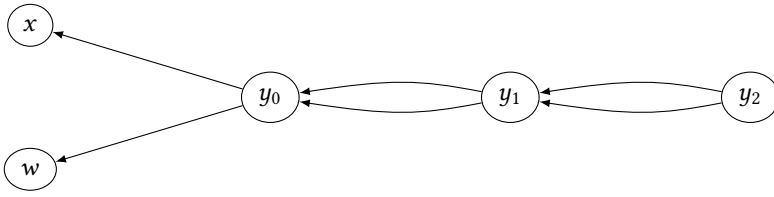


Fig. 3. A diamond dependent computational graph

We introduced a reference counting algorithm for dual numbers, to avoid the exponential time complexity, by only calling backward once.

The reference counting is managed in a wrapper of Tape, which has additional acquire and release functions.

Each wrapper has two internal states: (1) reference counter, (2) accumulator of delta. Respectively, acquire and release calls will increase and decrease the reference counter, and backward calls will cumulate the delta to the accumulator.

When the reference counting algorithm is enabled, backward is not recursive any more. Instead, a wrapper only call backward of its dependencies when the reference counter is decreased to zero. The entire process of backpropagation is shown in Figure 5.

This wrapper is implemented as the monadic data type Do, in which the side effects of updating counters and accumulators are monadic control flows. With the help of Do, now our computational graphs are modeled in Do[Tape[Data, Delta]], which can be created by forward methods described in Section 3.2. As mentioned in Section 3.3, computational graph node of binary operations are evaluated in parallel.

In a traditional backpropagation implementation, tape is a list, hence both the execution order of backward pass and forward pass must be sequentially reverse to each other. Even a previous attempt of closure-based tape [Pearlmutter and Siskind 2008] still requires conversion to sequential expressions of A-normal form [Sabry and Felleisen 1993].

By introducing reference counting, the execution order of our backward pass and forward pass do not have to be exactly reverse, hence the conversion to A-normal form becomes unnecessary. As a result, DeepLearning.scala supports out-of-order execution in both forward pass and backward pass, in which the independent sub-graph can be even executed in parallel.

5.5 The Overview of a Training Iteration

In brief, in each iteration, a differentiable function that contains multiple trainable variables can be trained in the following steps:

- (1) Executing the user-defined differentiable function with a batch of input, to obey a differentiable expression (i.e. a subtype of Layer).
- (2) Calling forward on differentiable expression to build a computational graph (i.e. a Do[Tape[Data, Delta]]). The reference counter to the computational graph is zero at the point.
- (3) Performing the forward pass of differentiable expression to build a tape (i.e. a Tape[Data, Delta]), which contains a pair of the result of forward pass and a backward closure. The reference counter of each node in a computational graph is increased during this step.
- (4) Performing backward closure of the root node of the computational graph. The accumulator of delta on the root node is updated.

834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882

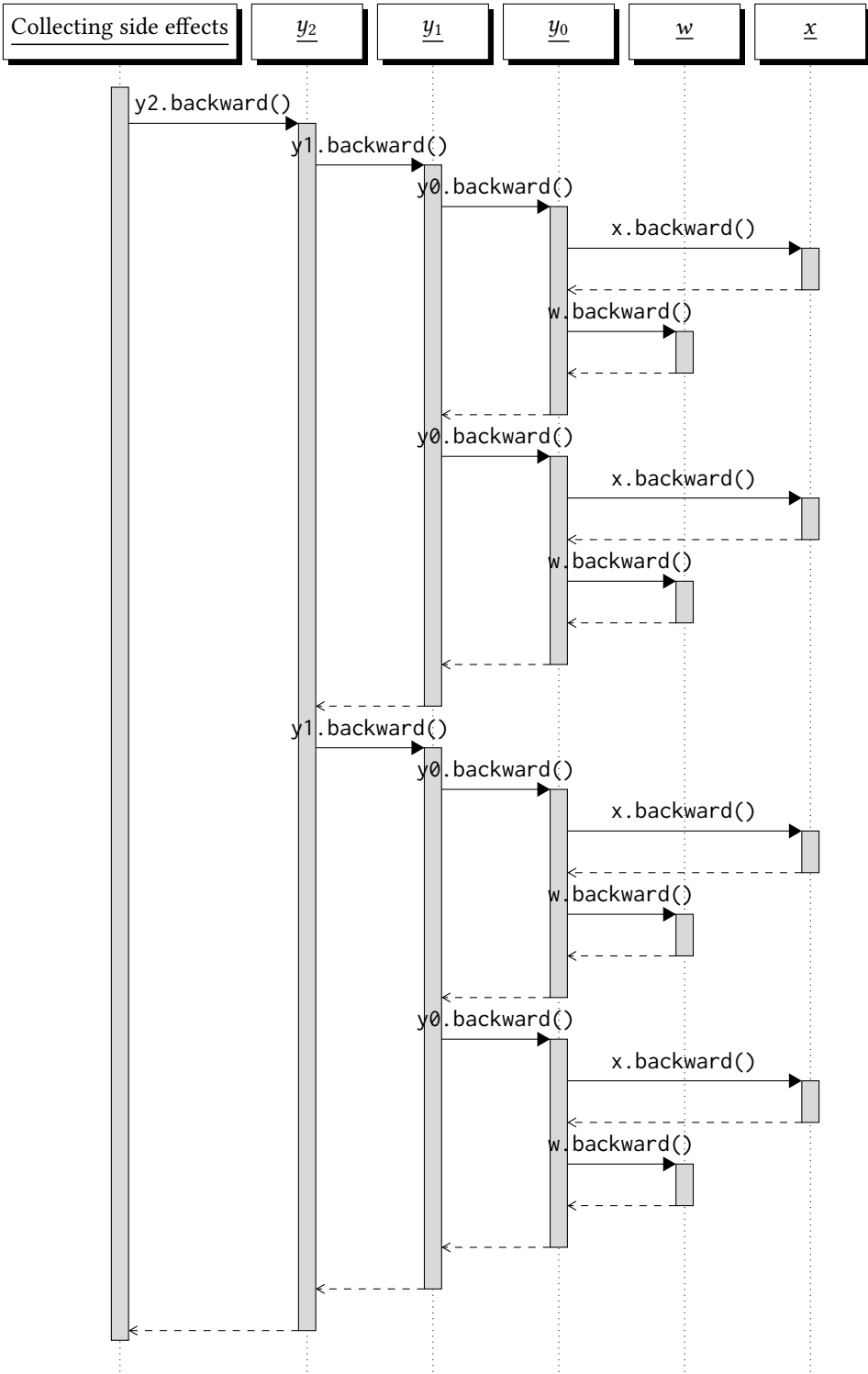


Fig. 4. Backpropagation for a diamond dependent computational graph

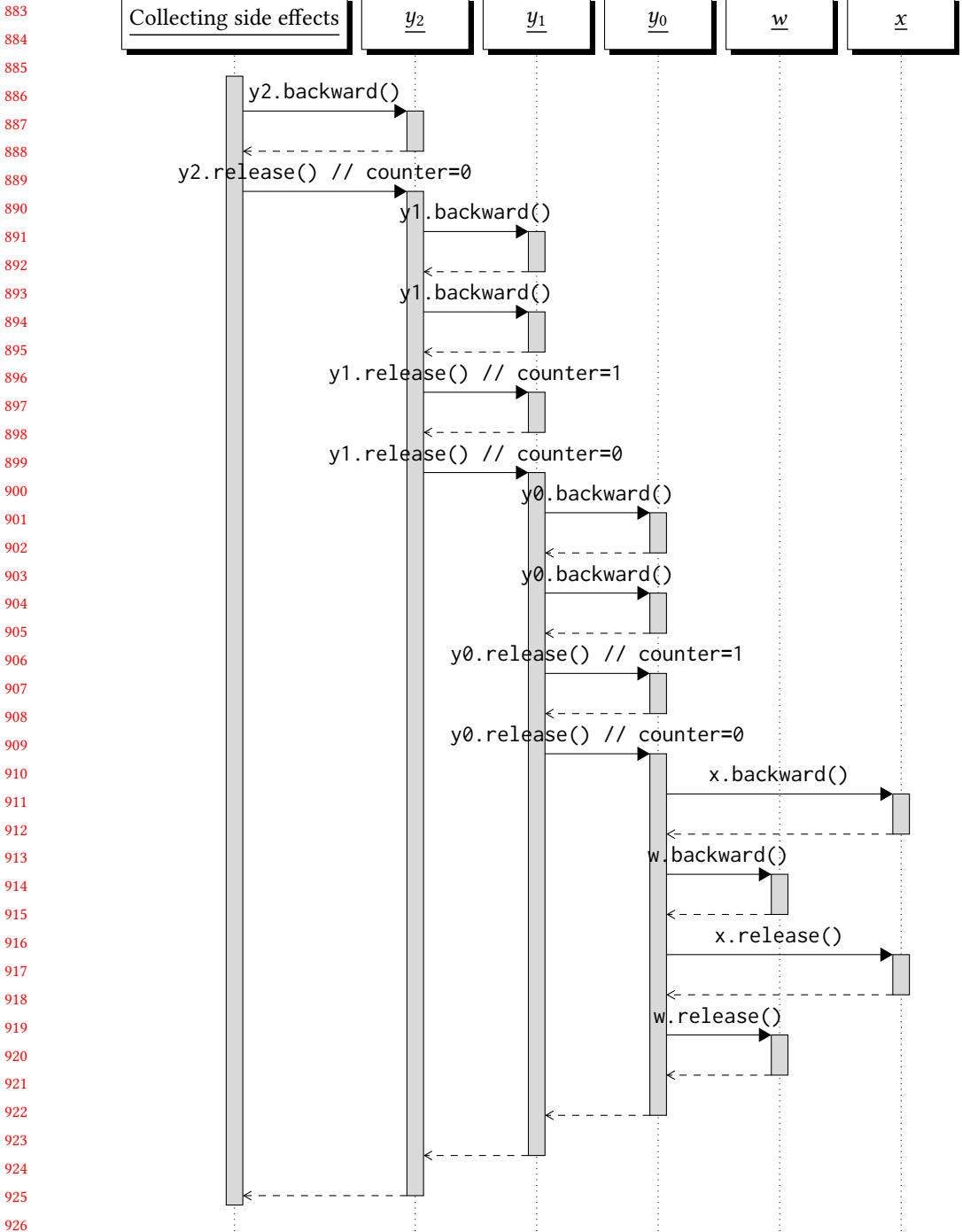


Fig. 5. Backpropagation for a diamond dependent computational graph (with reference counting)

- (5) Releasing of the root node of the **computational graph**. The reference counter of each node in a **computational graph** is decreased to zero and the backward closure of each node is performed during this step, thus all referenced **trainable variables** are updated.

Note that step 1 and step 2 are pure function calls, with no side effects. Step 3 to step 5 are monadic control flows, which encapsulate some side effects that will be performed only when an unsafe method `blockingAwait` or `onComplete` is eventually called.

6 EVALUATION

We created some benchmarks to evaluate the computational performance of `DeepLearning.scala`, especially, we want to measure:

- (1) How parallel execution affect the time cost of a training iteration of different structured neural networks;
- (2) the computational performance impact when a part of a neural network is disabled.

Those benchmarks are built with `Jmh`[Shipilev 2018], running on CPU, measuring the number of mini-batches per second, for training neural networks that contain variant number of branches of sub-networks, as classifiers for CIFAR-100[Krizhevsky and Hinton 2009], which is a dataset for image recognition, which has 100 fine-grained classes containing 600 images each. These fine-grained classes are grouped into 20 coarse-grained classes.

The architecture of the network used for benchmarks summarized in Figure 6.

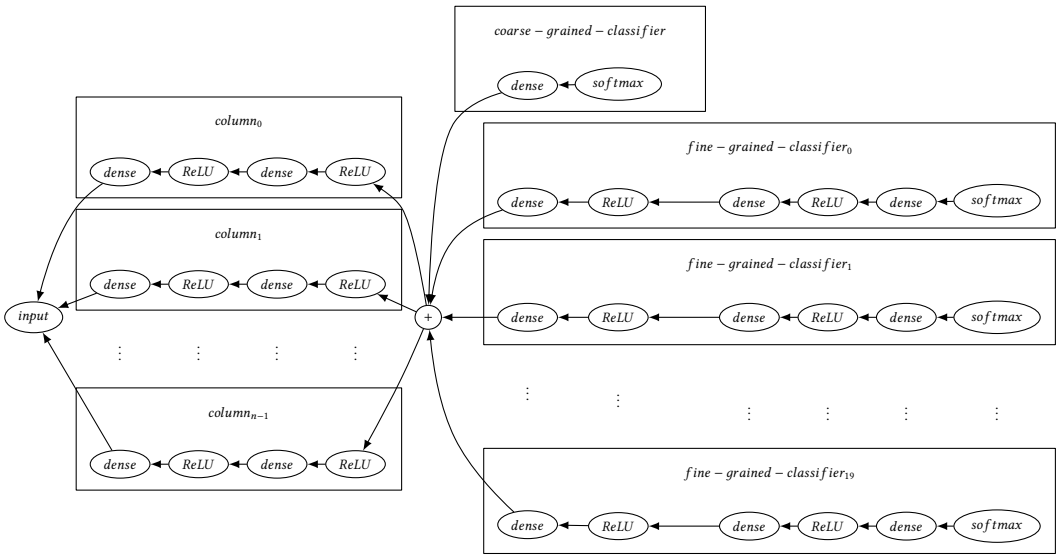


Fig. 6. The models used for benchmarks

We created n columns[Ciregan et al. 2012] of experts sub-networks for feature extracting. The output from those experts are then summed for further layers and classifiers. Each column contains two dense layers⁷ followed by ReLU activation layers. Those columns are independent, and we expect our framework can execute them in parallel when training or inference.

⁷ We use dense layers instead of convolution layers because the underlying N-dimensional array library ND4J is not able to efficiently perform immutable operation of convolution. See section 7.1 for discussion.

Since CIFAR-100 dataset has both coarse-grained and fine-grained labels, we created a coarse-grained classifier and 20 fine-grained classifiers. The coarse-grained classifier contains a dense layer to convert features to scores of 20 coarse-grained classes, followed by a softmax layer. Each fine-grained classifier corresponding to a coarse-grained class contains three dense layers, two ReLU activation layers, and a softmax layer, classifying 5 fine-grained classes.

We constructed mini-batches by a coarse-grained class when training. All samples in a mini-batch belongs to one coarse-grained class. With the help of the feature of dynamic neural network in DeepLearning.scala, for each training iteration, only one fine-grained classifier is used, other 19 fine-grained classifiers are skipped. When inferencing, the fine-grained classifier is chosen dynamically according to the prediction by the coarse-grained classifier. For comparison, we also created another set of benchmarks that do not skip unmatched fine-grained classifiers.

The benchmark result when training these models is shown in Table 2 and Table 3 (larger score is better). The sizes of those dense layers previous to softmax layers in these models are the number of classes (20 for coarse-grained classifier, 5 for fine-grained classifiers). Other dense layers output 64 features. These models are trained in a SGD optimizer with batch size 16. Both the number of columns in each model and the number of threads used in training vary among benchmarks.

number of columns	thread pool size	Score, ops/s
4	1	2.550 ± 0.043
4	2	3.233 ± 0.276
4	4	3.338 ± 0.217
2	1	3.345 ± 0.095
2	2	3.987 ± 0.091
2	4	4.828 ± 0.200
1	1	3.854 ± 0.046
1	2	5.239 ± 0.288
1	4	6.328 ± 0.058

Table 2. The benchmark result when no fine-grained classifier is skipped

number of columns	thread pool size	Score, ops/s
4	1	5.314 ± 0.304
4	2	5.484 ± 0.362
4	4	5.607 ± 0.159
2	1	8.702 ± 0.430
2	2	9.527 ± 0.259
2	4	8.831 ± 0.268
1	1	12.430 ± 1.514
1	2	13.531 ± 0.618
1	4	14.513 ± 0.554

Table 3. The benchmark result when unmatched fine-grained classifiers are skipped

The benchmark result verified the performance improvement when increasing thread pool size, since DeepLearning.scala executes independent sub-networks in parallel. This benchmark also shows a performance improvement when unmatched fine-grained classifiers are skipped.

7 FUTURE WORKS

7.1 New Back-end

Currently, DeepLearning.scala 2's built-in differentiable vector expression type `INDArrayLayer` is based on `nd4j`'s `INDArray`[Skymind 2017b]. As described in Section 5.5, in each training iteration, for each **computational graph** node, forward and backward operations are performed, which internally call some methods on `INDArray`, resulting in GPU kernel executions for `nd4j`'s CUDA runtime. These ND4J operations have bad computational performance because: (1) ND4J is not designed for immutable tensors;⁸ (2) some operations⁹ are extremely slow; (3) enqueueing a kernel is relatively expensive.

We are developing a new back-end as an alternative to `nd4j`. The new back-end will be able to merge multiple primitive operations into one larger kernel by dynamically generating OpenCL code. The new back-end will support more optimized operations on the GPU and reduce the number of kernel executions. We expect our new version will achieve better computational efficiency.

7.2 Distributed Model

Current DeepLearning.scala is only able to run on a standalone JVM, not a distributed cluster, thus it does not support "outrageously large neural networks"[Shazeer et al. 2017] that do not fit into the memory of a single node.

Since our **computational graph** is made of monadic expressions that consist of closures, they can be serialized and executed remotely in theory. We are investigating how to build a distributed machine learning system based on remotely executed monadic expression. We will find out if this suggested approach can support more complex models than the parameter server approach can.

8 DISCUSSION

DeepLearning.scala is an unique library among all deep learning frameworks. Our approach of AD has some attributes that never appear in previous frameworks.

8.1 Interoperable Differentiable Computational Graph

There were two different mechanisms in state-of-the-art deep learning frameworks: Define-and-Run vs. Define-by-Run.

State-of-the-art Define-and-Run frameworks[Abadi et al. 2016; Bergstra et al. 2010; Chen et al. 2015; Collobert et al. 2008; Intel 2016; Jia et al. 2014; Skymind 2017a] allows users to create **computational graphs**, which are immutable Abstract Syntax Trees (ASTs) of some object languages which can be evaluated by the framework runtime. Define-and-Run frameworks can schedule **computational graphs** to multiple CPUs, GPUs or other devices. However, the object languages have bad interoperability with the metalanguage. For example, a DeepLearning4j user cannot use Java control flows nor call Java native methods in neural networks.

State-of-the-art Define-by-Run frameworks[Google Brain 2017; Neubig et al. 2017; Paszke et al. 2017; Tokui et al. 2015] can eagerly execute actual forward pass calculation in user written code, and, at the same time, generate the internal states for running backward pass. Unlike Define-and-Run frameworks, Define-by-Run frameworks have good interoperability with the hosting language. Control flows and native function calls can be easily used during the execution of neural networks. However, Define-and-Run frameworks tend to store states and perform side effects when defining neural network structures, which makes this mechanism unsuitable for implementation in a purely functional flavor.

⁸<https://github.com/deeplearning4j/nd4j/issues/2271#issuecomment-344865418>

⁹`INDArray.broadcast` for example

1079 We discovered the third mechanism of monadic deep learning. Neural networks in DeepLearn-
 1080 ing.scala are immutable like in Define-and-Run frameworks, and interoperable with Scala like in
 1081 Define-by-Run frameworks.

1082 8.2 AD in a Functional Library

1084 Reverse mode AD as a functional library was previously considered as impossible to implement with-
 1085 out the ability to reflectively access and transform expressions associated with closures [Pearlmut-
 1086 ter and Siskind 2008]. For example, if you want to create a transform function that returns the
 1087 derivative for given function f :

```
1088 def transform(f: Double => Double): Double => Double
```

Listing 29. Impossible transform function for AD

1092 Obviously this transform function is impossible without the knowledge of the implementation
 1093 of f .

1095 Fortunately, in a statically typed language, the differentiable types and non-differentiable types
 1096 should differ for type safety. The type signature of our AD function can use the differentiable type
 1097 DoubleLayer instead of Double. It can be written as Listing 30:

```
1098 def typeSafeTransform(f: Double => DoubleLayer): Double => SideEffects = {  

  1099   input: Double =>  

  1100   val tape = f(input).forward  

  1101   tape.backward(Do.now(1.0))  

  1102 }  

  1103 }
```

Listing 30. Type safe transform function for AD

1104 Unlike [Pearlmutter and Siskind 2008]’s compiler primitive \overleftarrow{J} , our typeSafeTransform can use
 1105 the additional methods on DoubleLayer. As a result, our typeSafeTransform can be implemented
 1106 without reflection, as an ordinary Scala function, instead of a compiler primitive or a macro.

1109 We also change the derivative type to an opaque monadic type SideEffects. Unlike numeric
 1110 derivative, SideEffects can contain derivatives for more than one trainable variables, although
 1111 this change exclude the ability of higher order numerical differentiation.

1113 8.3 Portability

1114 Our approach can be implemented in a library, which requires few advanced language features.
 1115 Therefore, it can be ported to other languages that are not so powerful.

1117 Our static DAGs are applicative, which can be built from normal function calls or overloaded
 1118 operators without compiler-time transformation, as described in section 3.3.

1119 Our direct style dynamic neural networks require only one uncommon language feature, which
 1120 is !-notation, or direct style monadic expression, as described in section 3.1. The usage of this
 1121 feature does not prevent porting our approach to other languages, because:

- 1122 • The usage of !-notation is optional, as users of DeepLearning.scala can instead use **for**
 1123 comprehension if they choose not to use !-notation.
- 1124 • Direct style monadic expressions are available in Idris (!-notation), Haskell (do-notation) and
 1125 F# (computational expressions). Our approach can be ported to those languages without a
 1126 doubt.

1127

- All callback functions of monadic data types, including `Do`, `Future` and `UnitContinuation` will never be re-evaluated by `DeepLearning.scala` more than once. Therefore, those monadic data types can be replaced to generators, coroutines or one-pass continuations, which has become a widely adopted feature in many mainstream languages, including `ECMAScript`, `TypeScript`, `Go`, `Python`, `C#`, `C++20`, `Lua`.
- For Java or other languages that do not support generators or coroutines, those `!`-notation can be replaced to manually written `flatMap` calls.

8.4 Related works

In this section, we will discuss other deep learning frameworks in statically typed languages.

Several AD libraries [Baydin et al. 2015a; Bischof et al. 1992; Griewank et al. 1996; Hascoët and Pascual 2013] written in Fortran, C++ or F# support AD via operator overloading or external preprocessor but do not scale to deep neural networks, either due to using forward mode or lacking the feature of multiple trainable variables.

Other deep learning frameworks in statically typed languages (including Scala binding of C++ frameworks) [Baydin and Pearlmutter 2016; Chen 2017; Intel 2016; Skymind 2017a; Zhao et al. 2017] do not support AD, instead, they only provide their high level `computational graph` APIs to compose predefined layers into neural networks. As a result, those frameworks do not have the ability to create fine-grained custom algorithms.

[Elliott 2018] presented an approach that unified compiler-time translation for AD. The approach requires a compiler plug-in that translates Haskell programs into categorical form, which gives the ability to hook into each function call. By providing different instances of type classes, those function calls can translated to differential functions of different modes of AD. Unfortunately, the compiler plug-in is only available in Haskell, thus the approach is unable to be implemented in other languages. In contrast, our approach requires no metaprogramming features, thus can be implemented in many mainstream languages, as explained in section 8.3.

[Wang et al. [n. d.]] discovered an approach to perform reverse mode AD in delimited continuations. Their approach is very simple: (1) In forward pass, a chain of delimited continuations is created by `shift` calls, which inject some hooks that contain side effects to update the derivatives. (2) In backward pass, hooks are executed in reverse order, triggered by `reset`. However, their approach directly perform side effects in hooks, and always executes both forward pass and backward pass sequentially. In contrast, side effects in our approach are encapsulated in tasks of monadic data types, and we also provide nondeterministic implementation of applicative operations on those task types, allowing for parallel execution of both the forward pass and backward pass.

9 CONCLUSION

`DeepLearning.scala` is the first library that achieves all the following goals without metaprogramming technique:

- static type safety
- purely functional interface
- reverse mode AD
- multiple `trainable variables`
- interoperable internal DSL
- dynamic neural network
- statically type checking

With the help of `DeepLearning.scala`, normal programmers are able to build complex neural networks from simple code. They still write code as usual, and the only difference is that the code

1177 written in DeepLearning.scala is differentiable, which contains **trainable variables** that learn the
1178 knowledge.

1179

1180 GLOSSARY

1181 **computational graph** is an asynchronous monadic data type that manages the life cycle of tapes,
1182 whose type is `Do[Tape[Data, Delta]]` . 5, 7, 11, 14–20, 22, 24, 25

1183

1184 **differentiable expression** is a composable expression that supports operator overloading, whose
1185 type is `DoubleLayer`, `FloatLayer`, `INDArrayLayer`, or other subtypes of `Layer`. After a
1186 differentiable expression is built, it can perform forward pass to create a differentiable
1187 **computational graphs**. . 3, 5, 6, 10, 16, 17, 25

1188 **differentiable function** is a Scala function that returns a **differentiable expression**. It may repre-
1189 sent a loss functions, a neural network or a subset of a neural network (e.g. a dense block in
1190 `DenseNet`[[Iandola et al. 2014](#)]) . 5, 10, 11, 15–17

1191 **trainable variable** is a scalar or vector weight in a model, whose type is `DoubleWeight`, `FloatWeight`
1192 , `INDArrayWeight`, or other subtypes of `Weight` . 1, 3, 5, 10–12, 14–17, 20, 23–25

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

REFERENCES

- 1226
1227 Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat,
1228 Geoffrey Irving, Michael Isard, et al. 2016. TensorFlow: A System for Large-Scale Machine Learning.. In *OSDI*, Vol. 16.
1229 265–283.
- 1230 Atilim Gunes Baydin and Barak A Pearlmutter. 2016. Hype: Compositional Machine Learning and Hyperparameter
1231 Optimization. (2016). <https://hypelib.github.io/Hype/>
- 1232 Atilim Gunes Baydin, Barak A Pearlmutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind. 2015b. Automatic
1233 differentiation in machine learning: a survey. *arXiv preprint arXiv:1502.05767* (2015).
- 1234 Atilim Gunes Baydin, Barak A Pearlmutter, and Jeffrey Mark Siskind. 2015a. Diffsharp: Automatic differentiation library.
1235 *arXiv preprint arXiv:1511.07727* (2015).
- 1236 James Bergstra, Olivier Breuleux, Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, Guillaume Desjardins, Joseph Turian,
1237 David Warde-Farley, and Yoshua Bengio. 2010. Theano: A CPU and GPU math compiler in Python. In *Proc. 9th Python in
1238 Science Conf.* 1–7.
- 1239 Christian Bischof, Alan Carle, George Corliss, Andreas Griewank, and Paul Hovland. 1992. ADIFOR—generating derivative
1240 codes from Fortran programs. *Scientific Programming* 1, 1 (1992), 11–29.
- 1241 Tongfei Chen. 2017. Typesafe Abstractions for Tensor Operations. *arXiv preprint arXiv:1710.06892* (2017).
- 1242 Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu, Chiyuan Zhang, and Zheng
1243 Zhang. 2015. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv
1244 preprint arXiv:1512.01274* (2015).
- 1245 François Chollet et al. 2015. Keras. (2015).
- 1246 Dan Ciregan, Ueli Meier, and Jürgen Schmidhuber. 2012. Multi-column deep neural networks for image classification. In
1247 *Computer vision and pattern recognition (CVPR), 2012 IEEE conference on*. IEEE, 3642–3649.
- 1248 Ronan Collobert, K Kavukcuoglu, and C Farabet. 2008. Torch. In *Workshop on Machine Learning Open Source Software, NIPS*,
1249 Vol. 76.
- 1250 John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic
1251 optimization. *Journal of Machine Learning Research* 12, Jul (2011), 2121–2159.
- 1252 Conal Elliott. 2018. The simple essence of automatic differentiation. *Proceedings of the ACM on Programming Languages* 2,
1253 ICFP (2018), 70.
- 1254 Dumitru Erhan, Y Bengio, Aaron Courville, and Pascal Vincent. 2009. Visualizing Higher-Layer Features of a Deep Network.
1255 (01 2009).
- 1256 Martin Fowler. 2010. *Domain-specific languages*. Pearson Education.
- 1257 Google Brain 2017. *Eager Execution: An imperative, define-by-run interface to TensorFlow*. Google Brain. [https://research.
1258 googleblog.com/2017/10/eager-execution-imperative-define-by.html](https://research.googleblog.com/2017/10/eager-execution-imperative-define-by.html)
- 1259 Andreas Griewank, David Juedes, and Jean Utke. 1996. Algorithm 755: ADOL-C: a package for the automatic differentiation
1260 of algorithms written in C/C++. *ACM Transactions on Mathematical Software (TOMS)* 22, 2 (1996), 131–167.
- 1261 Dave Gurnell. 2017. *The Type Astronaut’s Guide to Shapeless*. Underscore Consulting LLP. [https://underscore.io/books/
1262 shapeless-guide/](https://underscore.io/books/shapeless-guide/)
- 1263 L. Hascoët and V. Pascual. 2013. The Tapenade Automatic Differentiation tool: Principles, Model, and Specification. *ACM
1264 Transactions On Mathematical Software* 39, 3 (2013). <http://dx.doi.org/10.1145/2450153.2450158>
- 1265 Forrest Iandola, Matt Moskewicz, Sergey Karayev, Ross Girshick, Trevor Darrell, and Kurt Keutzer. 2014. Densenet:
1266 Implementing efficient convnet descriptor pyramids. *arXiv preprint arXiv:1404.1869* (2014).
- 1267 Intel. 2016. BigDL. (2016). <https://github.com/intel-analytics/BigDL>
- 1268 Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor
1269 Darrell. 2014. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22nd ACM international
1270 conference on Multimedia*. ACM, 675–678.
- 1271 Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- 1272 Alex Krizhevsky and Geoffrey Hinton. 2009. Learning multiple layers of features from tiny images. (2009).
- 1273 Lanlan Liu and Jia Deng. 2017. Dynamic deep neural networks: Optimizing accuracy-efficiency trade-offs by selective
1274 execution. *arXiv preprint arXiv:1701.00299* (2017).
- 1275 Conor McBride and Ross Paterson. 2008. Applicative programming with effects. *Journal of functional programming* 18, 1
(2008), 1–13.
- 1276 Graham Neubig, Chris Dyer, Yoav Goldberg, Austin Matthews, Waleed Ammar, Antonios Anastasopoulos, Miguel Ballesteros,
1277 David Chiang, Daniel Clothiaux, Trevor Cohn, et al. 2017. DyNet: The Dynamic Neural Network Toolkit. *arXiv preprint
1278 arXiv:1701.03980* (2017).
- 1279 BCDS Oliveira, A Moors, and M Odersky. 2010. Type classes as objects and implicits. *ACM SIGPLAN Notices* (2010).
- 1280 Erik Osheim and Jorge Vicente Cantero. 2017. *SIP-35 - Opaque types*. <https://docs.scala-lang.org/sips/opaque-types.html>

- 1275 Adam Paszke, Sam Gross, Soumith Chintala, and Gregory Chanan. 2017. *PyTorch: Tensors and Dynamic neural networks in*
1276 *Python with strong GPU acceleration*. <http://pytorch.org/>
- 1277 Barak A Pearlmutter and Jeffrey Mark Siskind. 2008. Reverse-mode AD in a functional framework: Lambda the ultimate
1278 backpropagator. *ACM Transactions on Programming Languages and Systems (TOPLAS)* 30, 2 (2008), 7.
- 1279 David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. 1985. *Learning internal representations by error propagation*.
Technical Report. California Univ San Diego La Jolla Inst for Cognitive Science.
- 1280 Amr Sabry and Matthias Felleisen. 1993. Reasoning about programs in continuation-passing style. *Lisp and symbolic*
1281 *computation* 6, 3-4 (1993), 289–360.
- 1282 Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outra-
1283 geously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538* (2017).
- 1284 Aleksey Shipilev. 2018. JMH: Java Microbenchmark Harness. (2018). <http://openjdk.java.net/projects/code-tools/jmh/>
1285 <http://openjdk.java.net/projects/code-tools/jmh/>
- 1286 SkyMind 2017a. *Deeplearning4j: Open-Source, Distributed, Deep Learning Library for the JVM*. SkyMind. <https://deeplearning4j.org/>
- 1287 SkyMind 2017b. *ND4J: N-Dimensional Arrays for Java*. SkyMind. <https://nd4j.org/>
- 1288 Seiya Tokui, Kenta Oono, Shohei Hido, and Justin Clayton. 2015. Chainer: a next-generation open source framework for
1289 deep learning. In *Proceedings of workshop on machine learning systems (LearningSys) in the twenty-ninth annual conference*
1290 *on neural information processing systems (NIPS)*, Vol. 5.
- 1291 Fei Wang, James Decker, Xilun Wu, Grégory Essertel, and Tiark Rompf. [n. d.]. Backpropagation with Continuation
1292 Callbacks: Foundations for Efficient and Expressive Differentiable Programming. ([n. d.]).
- 1293 Bo Yang. 2017. *Dsl.scala: a framework to create DSL in Scala*. (2017). <https://github.com/ThoughtWorksInc/Dsl.scala/>
- 1294 Kenji Yoshida. 2017. *Scalaz: An extension to the core scala library*. <https://scalaz.github.io/scalaz/>
- 1295 Matthew D Zeiler. 2012. ADADELTA: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701* (2012).
- 1296 Tian Zhao, Xiaobing Huang, and Yu Cao. 2017. DeepDSL: A Compilation-based Domain-Specific Language for Deep
1297 Learning. *arXiv preprint arXiv:1701.02284* (2017).
- 1298
- 1299
- 1300
- 1301
- 1302
- 1303
- 1304
- 1305
- 1306
- 1307
- 1308
- 1309
- 1310
- 1311
- 1312
- 1313
- 1314
- 1315
- 1316
- 1317
- 1318
- 1319
- 1320
- 1321
- 1322
- 1323