

# Correctness is Demanding, Performance is Frustrating

ANONYMOUS AUTHOR(S)

In this paper we demonstrate a technique for developing high performance applications with strong correctness guarantees. We use a theorem prover to derive a high-level specification of the application that includes correctness invariants of our choice. After that, within the same theorem prover, we implement an extraction of the specified application into a high-performance language of our choice. Concretely, we are using Agda to specify a framework for automatic differentiation (reverse mode) that is focused on index-safe tensors. This framework comes with an optimiser for tensor expressions and the ability to translate these expressions into SaC and C. We specify a canonical convolutional neural network within the proposed framework, compute the derivatives needed for the training phase and then demonstrate that the generated code matches the performance of hand-written code when running on a multi-core machine.

Additional Key Words and Phrases: Dependent Types, Agda, Array Programming, Automatic Differentiation, SaC

## ACM Reference Format:

Anonymous Author(s). 2025. Correctness is Demanding, Performance is Frustrating. 1, 1 (February 2025), 25 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

## 1 Introduction

The year is 2024, and we still have to make a formidable choice between correctness and performance for all the programming projects that we start. Low level programming languages such as C or Fortran make it possible to leverage intricate hardware features for the price of poor analysability and correctness guarantees. Dependently-typed systems such as Lean or Agda make it possible to describe arbitrary invariants within the given program, yet they can rarely generate high-performance code.

The idealist inside of us is exasperated, because there should be a perfect solution that caters for both cases. However, the practitioner within us proposes a different perspective. Instead of using a single language, we may use two languages in cooperation. Specifically, we envision an alliance between a proof assistant and some high-performance language of our choice.

To explore this idea, we investigate the concrete problem of automatic differentiation (AD) which is often found in machine learning applications. This is a convenient case study as it comes with the following challenges. From the correctness perspective, it is crucially important to track the shapes and ranks of the tensors, guaranteeing the absence of out-of-bound indexing. This is a very common source of errors that can be incredibly difficult to find. Secondly, we have to compute derivatives of the given tensor expressions, preserve safe indexing guarantees while we do so, and we have to be able to translate the computed expressions into some high-performance language. As machine learning applications are known to be numerically intensive problems, our performance challenge lies in running the program as fast as we can on the chosen hardware architecture.

We follow [23] which demonstrates that it is possible to implement one of the canonical convolutional neural network (CNN) in the array language SaC [9, 20], obtaining good sequential and

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM XXXX-XXXX/2025/2-ART

<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

parallel performance that is competitive with TensorFlow [1] and PyTorch [16]. Focussing on correctness, we propose a theory of rank-polymorphic arrays [21] in Agda [2]. Within this framework, we encode the CNN from [23] and lift it into an embedded DSL. We implement AD (reverse mode) and domain-specific optimisations for expressions in that DSL. Finally, we implement an extraction into SaC (functional array language) and C (low-level imperative language).

As a result, we demonstrate an approach where the entire specification, optimiser, AD and code generation are available to us within a proof assistant of our choice. We can prove facts about all the stages of the pipeline and easily adjust them to our liking. We argue that such a liberating approach is feasible in practice, at least for the times of dialectic of correctness and performance.

The contributions of this paper are as follows:

- (1) a rank-polymorphic array theory in Agda;
- (2) an implementation of the CNN from [23] in Agda;
- (3) an embedded DSL in Agda which supports AD (reverse mode);
- (4) an extraction mechanism for generating SaC or C code from the DSL; and
- (5) an experimental evaluation of the generated codes.

This paper is written in literate Agda, which guarantees that all the code snippets have been type-checked.

## 2 Background

Automatic differentiation has been around for many decades [10, 15], so it is well-understood at a conceptual level. However, a number of questions related to bringing AD into the context of programming languages remain open. Recent successes in machine learning have spurred further interest in AD which has led to several new developments. For the context of this paper, we focus on recent work that contributes to the perspective of balancing correctness guarantees and performance. Our selection here is by no means exhaustive, for a broader scope we refer the reader to [3].

There has been a number of programming-language-oriented approaches that explain how to add AD to a programming language of choice. Examples of these include Futhark [18], Haskell [12], and Jax [13]. Furthermore, a number of machine learning frameworks that incorporate AD have been proposed in recent years: TensorFlow [1], PyTorch [16], MXNet [5] and many more. While in particular the dedicated frameworks tend to find widespread acceptance by practitioners, both, correctness and performance leave two open questions: (i) how is it possible to ensure that the AD algorithm is implemented correctly? (ii) if the language or the framework do not perform as expected, what are the options to solve this? Unfortunately, for many cases the answer to both questions is unsatisfying. Most of the languages/frameworks do not come with formal correctness guarantees, so one has to trust the implementers of these tools. One can run tests as well to gain trust in the implementation but that is far from a formal guarantee. If one relies on the AD provided by a chosen language/framework, and the generated code does not perform well, one has to modify the language/framework, as these solutions are tightly integrated with the tools. The problem here is that most of these tools have very large and sophisticated implementations typically comprising of hundreds of thousands of lines of code. Furthermore, these systems often rely on sophisticated combinations of sub-systems that need to be fine-tuned to the executing hardware.

Another line of work studies high-level principles of AD using category theory [6, 8, 22]. While this indeed comes with great correctness guarantees due to some naturality principles, it is not always clear how to implement this in a way that leads to efficiently executable specifications. Also, the entire treatment of index-safe tensors is unclear.

In [7] the author proposes to view AD problem using the language of cartesian categories. It has been shown that this approach can be used in practice by implementing the proposed technique in Haskell. Type classes are a vehicle to restrict expressions that are differentiable. There is a Haskell plugin that translates expressions that are instances of the mentioned type classes into categorical primitives, AD is performed on these and the code is reflected back to Haskell. This is a nice approach that makes it easy to verify the correctness of the algorithm. However, the treatment of tensors and general extractability remains a little unclear. While it is briefly mentioned that representable functors are supported, it is unclear whether this is sufficient to represent rank-polymorphic arrays with strict bound checks. Also, correctness guarantees are inevitably restricted by the Haskell type system, so we are likely to find invariants that are inexpressible in that setup.

### 3 Array Theory

The central data structures of our case study are multi-dimensional arrays. This section is dedicated to defining a minimalist array theory in Agda which is well-suited for a specification of CNNs.

We assume that the reader is sufficiently familiar with Agda's syntax. For gentle introductions we refer to one of the tutorials that are freely available online<sup>1</sup>.

The conciseness of the specification in [23] relies on rank-polymorphism, which is the ability to operate on arrays of arbitrary rank. We capture this feature in our array theory. The central consideration when working with dependent types is how to represent data. Some encodings are better suited for reasoning, others are more efficient at runtime. Due to our two-language setup, our choice of representation is driven by proof considerations only. This enables us to represent arrays as functions from indices to values.

An absence of out-of-bound errors requires that all array indices fall within the shapes of the arrays that they are selecting from. The shape of an array describes the extent of each of its axes. We represent shapes lists of natural numbers using the data type `S`. Empty shapes `[]` correspond to singleton arrays (often called scalars in array languages). The cons operation `_::_` prepends the axis to the left of the shape. Note that underscores in `_::_` specify the position where arguments go, therefore `::` is an infix binary operation.

Array positions (indices) are given by the dependent type `P` which is indexed by shapes. A position within an array of shape `s` has exactly the same tree structure as `s`, but the leaves are natural numbers that are bounded by the corresponding shape elements.

Arrays are given by the data type `Ar` which is indexed by a shape and an element type. The formal definitions of `S`, `P` and `Ar` are as follows:

```
data S : Set where
  [] : S
  _::_ : ℕ → S → S

data P : S → Set where
  [] : P []
  _::_ : Fin n → P s → P (n :: s)

Ar : S → Set → Set
Ar s X = P s → X
```

As arrays are functions, selections are function applications and array construction is a function definition (e.g. via  $\lambda$ -abstraction).

*Array Combinators.* It is helpful to invest a little time in defining array combinators. First, we can observe that `Ar` of a fixed shape is an applicative functor [14], so we can trivially derive: `K x` to produce a constant array; `map f a` to apply `f` to all the elements of `a`; and `zipWith f a b` to

<sup>1</sup>See <https://agda.readthedocs.io/en/v2.6.3/getting-started/tutorial-list.html>.

point-wise apply the binary operation  $f$  to  $a$  and  $b$ .

$$\begin{aligned} \text{K} : X \rightarrow \text{Ar } s X & & \text{map} : (X \rightarrow Y) \rightarrow \text{Ar } s X \rightarrow \text{Ar } s Y \\ \text{K } x \ i = x & & \text{map } f \ a \ i = f \ (a \ i) \\ \\ \text{zipWith} : (X \rightarrow Y \rightarrow Z) \rightarrow \text{Ar } s X \rightarrow \text{Ar } s Y \rightarrow \text{Ar } s Z \\ \text{zipWith } f \ a \ b \ i = f \ (a \ i) \ (b \ i) \end{aligned}$$

Array shapes can be concatenated as lists. We call this operation *shape product* and we denote it with  $\_ \otimes \_$  (because this correspond to the shape of tensor product). Positions of sub-shapes can be joined into a position of a product shape using the  $\_ \otimes_p \_$  operation. Dually, positions of a product shape can be split into positions of the corresponding subshapes using  $\text{split}$ . Here are the types of these three operations:

$$\_ \otimes \_ : S \rightarrow S \rightarrow S \quad \_ \otimes_p \_ : P \ s \rightarrow P \ p \rightarrow P \ (s \otimes p) \quad \text{split} : P \ (s \otimes p) \rightarrow P \ s \times P \ p$$

Arrays are homogeneously nested, *i.e.* the shapes of all the sub-arrays have to be the same. Therefore, we can switch between the array of a product shape and the nested array (array of arrays). This operation is very similar to currying except it happens at the level of shapes. The combinators that achieve this are named  $\text{nest}$  and  $\text{unnest}$  and their definitions are:

$$\begin{aligned} \text{nest} : \text{Ar } (s \otimes p) \ X \rightarrow \text{Ar } s \ (\text{Ar } p \ X) & & \text{unnest} : \text{Ar } s \ (\text{Ar } p \ X) \rightarrow \text{Ar } (s \otimes p) \ X \\ \text{nest } a \ i \ j = a \ (i \otimes_p j) & & \text{unnest } a \ i = \text{uncurry } a \ (\text{split } i) \end{aligned}$$

*Reduction.* We implement reduction of the binary operations over arrays in two steps. Firstly, we define 1-d reductions that we call  $\text{sum}_1$  which is very similar to right fold on lists. The arrays of higher ranks iterate  $\text{sum}_1$  bottom-up. The definition of the primitives are as follows:

$$\begin{aligned} \text{pattern } \iota \ n = n :: [] & & \text{sum}_1 : (X \rightarrow X \rightarrow X) \rightarrow X \rightarrow \text{Ar } (\iota \ n) \ X \rightarrow X \\ \iota \text{suc} : P \ (\iota \ n) \rightarrow P \ (\iota \ (\text{suc } n)) & & \text{sum}_1 \ \{n = \text{zero}\} \ f \ \epsilon \ a = \epsilon \\ \iota \text{suc} \ (\iota \ i) = \iota \ (\text{suc } i) & & \text{sum}_1 \ \{n = \text{suc } n\} \ f \ \epsilon \ a = f \ (a \ (\iota \ \text{zero})) \ (\text{sum}_1 \ f \ \epsilon \ (a \circ \iota \text{suc})) \\ \\ \text{sum} : (X \rightarrow X \rightarrow X) \rightarrow X \rightarrow \text{Ar } s \ X \rightarrow X \\ \text{sum} \ \{s = []\} \ f \ \epsilon \ a = f \ \epsilon \ (a \ []) \\ \text{sum} \ \{s = x :: s\} \ f \ \epsilon \ a = \text{sum}_1 \ f \ \epsilon \ \$ \ \text{map} \ (\text{sum } f \ \epsilon) \ (\text{nest } a) \end{aligned}$$

Note that our reduction forces the types of the arguments of the binary operation to be the same, which is different from the usual  $\text{foldr}$  definition. While we do not need this functionality for our example, it is worth noting that the standard behaviour can be recovered<sup>2</sup> through reduction of function composition.

## 4 CNN Building Blocks

With the array theory from the previous section we can define the actual primitives that are required for our case study.

<sup>2</sup>We recover regular fold behaviour by running  $\text{sum}$  over function composition:

$$\begin{aligned} \text{sum}' : (X \rightarrow Y \rightarrow Y) \rightarrow Y \rightarrow \text{Ar } s \ X \rightarrow Y \\ \text{sum}' \ f \ \epsilon \ a = \text{sum } \_ \circ \_ \text{id} \ (\text{map } f \ a) \ \epsilon \end{aligned}$$

## 4.1 One-dimensional convolution

We start with plus and minus operations for 1-d indices which will be used in the definition of convolution:

$$\begin{aligned} \_ \oplus \_ : \text{Fin } m \rightarrow \text{Fin } (1 + n) \rightarrow \text{Fin } (m + n) & \quad \_ \ominus \_ : (i : \text{Fin } (m + n)) (j : \text{Fin } m) \\ \text{zero} \oplus j = \text{inject-left } j & \quad \rightarrow \text{Dec } (\exists \lambda k \rightarrow j \oplus k \equiv i) \\ \text{suc } i \oplus j = \text{suc } (i \oplus j) & \end{aligned}$$

While the definitions look very innocent, their types carry non-trivial information. Consider  $\_ \oplus \_$  which is addition of bounded  $i$  and  $j$ . However, the type says:

$$\frac{i < m \quad j < 1 + n}{i + j < m + n}$$

This looks a little surprising, but this indeed holds for natural numbers. A reader may convince herself by considering the maximum value that  $i$  and  $j$  can possibly take. The  $\_ \oplus \_$  have partial inverses making it possible to define left and right subtraction. We consider left subtraction  $\_ \ominus \_$ . Its type says that there exists a decision procedure for finding  $k$  of type  $\text{Fin } (1 + n)$  together with the proof that  $k$  is an inverse. In some sense  $\text{Dec}$  is similar to  $\text{Maybe}$  type, except it forces one to prove why the value does not exist as opposed to just returning  $\text{nothing}$ . This happens to be very useful, as it is really easy to introduce off-by-one errors otherwise.

With the above definitions we can define convolution for 1-dimensional cases. A side note for mathematically inclined readers. We use the term *convolution* in the way it is used in machine learning. Technically, we compute a cross-correlation, because the array of weights is not flipped. However, in practice this is not a problem, as we assume that weights are stored flipped in memory.

We define a handy shortcut  $\text{Vec}$  and  $\text{Ix}$  which are  $\text{Ar}$  and  $\text{P}$  for 1-dimensional cases.

$$\begin{aligned} \text{Vec} : \mathbb{N} \rightarrow \text{Set} \rightarrow \text{Set} & \quad \text{Ix} : \mathbb{N} \rightarrow \text{Set} \\ \text{Vec } m \ X = \text{Ar } (\iota \ m) \ X & \quad \text{Ix } m = \text{P } (\iota \ m) \end{aligned}$$

We introduce the  $\text{slide}_1$  primitive that selects a  $(1 + n)$ -element vector from the  $(m + n)$ -element vector starting at the offset  $i$ . Then, following [23], we compute  $m$ -element array of slides and then sum it up.

$$\begin{aligned} \text{slide}_1 : \text{Ix } m \rightarrow \text{Vec } (m + n) \ X \rightarrow \text{Vec } (1 + n) \ X \\ \text{slide}_1 (\iota \ i) \ v (\iota \ j) = v (\iota \ (i \oplus j)) \\ \text{conv}_1 : \text{Vec } (m + n) \ \mathbb{N} \rightarrow \text{Vec } m \ \mathbb{N} \rightarrow \text{Vec } (1 + n) \ \mathbb{N} \\ \text{conv}_1 \ a \ w = \text{sum } (\text{zipWith } \_ \_ ) (\text{K } 0) (\lambda \ i \rightarrow \text{map } (w \ i \ * \ \_) (\text{slide}_1 \ i \ a)) \end{aligned}$$

## 4.2 Generalisation

We want to define convolution for arrays of higher ranks. The first thing to do is to express  $m + n$  and  $1 + n$  where  $m$  and  $n$  become arbitrary shapes. In case of addition, we need a witness that both shapes have the same length. If they do, we can simply add their nodes point-wise. We define the three-way relation  $\_ + \_ \approx \_$  that combines the witness and the action. That is, the type  $p + q \approx r$  says that  $p$  and  $q$  have the same length and that  $q$  is a point-wise addition of  $p$  and  $q$ . We introduce a similar relation  $\text{suc } \_ \approx \_$  for  $1 + n$  case, and we introduce the relation  $\_ * \_ \approx \_$  that witnesses point-wise multiplication that will be needed for blocking. We define these relation in two steps. Firstly, we give a generalised pointwise-relation for binary and ternary relations on natural numbers:

**data**  $\text{Pointw}_2 \ (R : (a \ b : \mathbb{N}) \rightarrow \text{Set}) : (a \ b : \mathbb{S}) \rightarrow \text{Set}$  **where**  
**instance**

```

246 [] : Pointw2 R [] []
247 cons : { R m n } → { Pointw2 R s p } → Pointw2 R (m :: s) (n :: p)
248

```

```

249 data Pointw3 (R : (a b c : ℕ) → Set) : (a b c : S) → Set where
250 instance

```

```

251 [] : Pointw3 R [] [] []
252 cons : { R m n k } → { Pointw3 R s p q } → Pointw3 R (m :: s) (n :: p) (k :: q)
253

```

While the definition is straight-forward, note that we mark constructors with the keyword `instance` and we turn the arguments of `cons` into instance arguments<sup>3</sup>. These arguments behave like the hidden arguments, except Agda will apply an instance search when solving them. This allows us to omit these proofs in a larger number of cases when compared to hidden arguments.

Definition of the actual relations are:

```

254 _+_≈_ : (s p q : S) → Set
255 s + p ≈ q = Pointw3 (λ x y z → x + y ≡ z) s p q
256

```

```

257 _*_≈_ : (s p q : S) → Set
258 s * p ≈ q = Pointw3 (λ x y z → x * y ≡ z) s p q
259
260 suc_≈_ : (s p : S) → Set
261 suc s ≈ p = Pointw2 (λ x y → suc x ≡ y) s p
262

```

With these relations in place, how do we define generalised convolution? One possible way is to use the `sum` approach where we recurse over the shape tree and perform one operation at a time. However, there is a good point made in [23] that we can shift the shape recursion into index computation. Therefore we define `_⊕p_` and `_⊖p_` which generalise `_⊕_` and `_⊖_` for higher ranks. Once again, `Dec` type forces `⊖p` to justify the cases when the inverse does not exist.

```

263 _⊕p_ : P s → P u → suc p ≈ u → s + p ≈ r → P r
264 (i ⊕p j) [] [] = j
265 ((i :: is) ⊕p (j :: js)) (cons { refl } { sp }) (cons { refl } { s+p })
266 = (i ⊕p j) :: (is ⊕p js) sp s+p
267
268 _⊖p_ : (i : P r) (j : P s) (su : suc p ≈ u) (sp : s + p ≈ r) → Dec (∃ λ k → (j ⊕p k) su sp ≡ i)
269

```

We do not show the implementation of the `⊖p`, but it very much follows the structure of `⊕p`: we apply `⊖` on the leaves and we recurse on the product shape with a little bit plumbing to construct the proof of (non-)existence of the inverse.

Our generalised `slide` looks very much the same as its 1-dimensional counterpart. All the difference lies in the index computation. We also introduce a section of the slide that we call `backslide` which embed a  $(1 + p)$ -dimensional array into a  $(s + p)$ -dimensional one at offset  $i$  using some the provided default element `def`.

```

270 slide : P s → s + p ≈ r → Ar r X → suc p ≈ u → Ar u X
271 slide i pl a su j = a ((i ⊕p j) su pl)
272
273 backslide : P s → Ar u X → suc p ≈ u → (def : X) → s + p ≈ r → Ar r X
274 backslide i a su def pl j with ((j ⊖p i) su pl)
275 ... | yes (k , _) = a k
276 ... | _ = def
277

```

<sup>3</sup>See <https://agda.readthedocs.io/en/v2.7.0.1/language/instance-arguments.html> for more details.

### 4.3 Remaining primitives

In the rest of this section we implement the remaining CNN-specific primitives. We are going to use the builtin Float type that we call  $\mathbb{R}$  so that we can run our specification with concrete values. However, all we require from  $\mathbb{R}$  is a set of standard arithmetic operations. Therefore,  $\mathbb{R}$  can be abstracted out as a parameter.

Generalised convolution is given by `conv`, and it is almost identical to its 1-dimensional counterpart (except it used `slide` instead of `slide1`). The `mconv` runs `u convs` adds biases to each of them from the array `b`.

```
conv : s + p ≈ r → Ar r ℝ → Ar s ℝ → suc p ≈ u → Ar u ℝ
conv sp a w su = sum (zipWith _+_ ) (K 0.0) λ i → map (w i *_) (slide i sp a su)
```

```
mconv : { { s + p ≈ r } } → Ar r ℝ → Ar (u ⊗ s) ℝ → Ar u ℝ → { { suc p ≈ q } } → Ar (u ⊗ q) ℝ
mconv { { sp } } inp w b { { su } } = unnest λ i → map (b i +_) (conv sp inp (nest w i) su)
```

The logistic function computes  $1/(1 + e^{-x})$  for every element in the array.

```
logistic : Ar s ℝ → Ar s ℝ
logistic = map λ x → 1.0 ÷ (1.0 + e^ (- x))
```

*Average Pooling.* One of the steps of the machine learning algorithm is average pooling which splits an array into sub-blocks and computes the average for every such block. Implementing this pattern generally is tricky as we have to preserve local neighbourhood within the blocks. Working with a pre-blocked array would be inconvenient as the blocked shaped does not go well with `slides`. We solvet this by introducing blocked selections into arrays of shape  $(s * p)$  as well and blocked array constructor `imapb` that builds an array of shape  $(s * p)$  out of `s` blocks of shape `p`. Defining these operations we require pairing and projections of the blocked indices:

```
ix-div : P q → s * p ≈ q → P s          ix-mod : P q → s * p ≈ q → P p
ix-combine : P s → P p → s * p ≈ q → P q
```

With these operations, definitions of `selb` and `imapb` are:

```
selb : Ar q X → p * s ≈ q → P p → Ar s X          imapb : Ar s (Ar p X) → s * p ≈ q → Ar q X
selb a p i j = a (ix-combine i j p)                  imapb a p i = a (ix-div i p) (ix-mod i p)
```

Finally we define an average pooling that is specialised to the 2-dimensional case, that is needed in our running example.

```
avgp2 : (m n : ℕ) → Ar (m ℕ.* 2 :: n ℕ.* 2 :: []) ℝ → Ar (m :: n :: []) ℝ
avgp2 m n a = map ((_ ÷ fromN 4) ∘ sum _+_ 0.0) (selb a it)
```

We are now ready to provide the implementation of the forward part of the CNN as follows. The `inp` argument is the image of a hand-written digit, all the other arguments are weights, and the



function returns the 10-element vector with probabilities which digit that is.

```

forward : (inp : Ar (28 :: 28 :: []) ℝ) → (k1 : Ar (6 :: 5 :: 5 :: []) ℝ)
        → (b1 : Ar (6 :: []) ℝ) → (k2 : Ar (12 :: 6 :: 5 :: 5 :: []) ℝ)
        → (b2 : Ar (12 :: []) ℝ) → (fc : Ar (10 :: 12 :: 1 :: 4 :: 4 :: []) ℝ)
        → (b : Ar (10 :: []) ℝ) → Ar (10 :: 1 :: 1 :: 1 :: 1 :: []) ℝ
forward inp k1 b1 k2 b2 fc b = let
  c1 : Ar (6 :: 24 :: 24 :: []) ℝ
  c1 = logistic $ mconv inp k1 b1

  s1 : Ar (6 :: 12 :: 12 :: []) ℝ
  s1 = unnest {s = 6 :: []} $ map (avgp2 12 12) (nest c1)

  c2 : Ar (12 :: 1 :: 8 :: 8 :: []) ℝ
  c2 = logistic $ mconv s1 k2 b2

  s2 : Ar (12 :: 1 :: 4 :: 4 :: []) ℝ
  s2 = unnest {s = 12 :: 1 :: []} $ map (avgp2 4 4) (nest c2)

  r = logistic $ mconv s2 fc b
in r

```

## 5 Embedded DSL

Any implementation of automatic differentiation has to decide which operations are supported. Surely, it does not make sense to compute derivatives of a function that opens a file. This choice, no matter how it is implemented, can be seen as a definition of an embedded language. Once we accept to identify an embedded language, the idea of embedding it in a way that facilitates extraction actually appears rather naturally and thus advances the approach that we propose in this paper.

Coming back to our example, we have to choose the primitives that the embedded language should support. They need to be sufficient to express AD as well as to define CNNs. The main trade-off here is the choice of the level of abstraction of these primitives: low-level primitives are easier to differentiate, but they make the overall expressions more complex which also adds to the challenge of optimising code. Making this choice is difficult and, most likely, requires quite some adjustment when striving for performance. Here we see a key benefit of the approach we propose in this paper: the use of a single framework for the embedding, the optimisation, and the extraction makes the implementation comparatively small, allowing for quick adjustments in the level of abstraction, code optimisation and its extraction. We start with a pragmatic approach; we include the primitives that are either shared by the model and the back-end or that can be easily implemented in the back-end language.

It turns out that it is possible to choose the primitives in a way that the derivatives can be expressed in the very same embedded language. While this at first glance may just seem to be just a nice coincidence, it turns out that this has several tangible benefits: high-order derivatives can be computed by the same transformation and we can share all optimisations between the code itself and its derivatives.

As we operate within a dependently-typed proof-assistant, we can easily make our embedded language well-scoped and intrinsically typed (shaped in our case). Our context `Ctx` is a snoc-list of shapes where each shape has a tag indicating whether it is an index or an array. We use de Bruijn variables given by the relation `_∈_` in the usual way. We also define variables  $v_1$ ,  $v_2$ , *etc.* by



iteratively applying  $v_s$  to  $v_0$  (definition not shown).

```

data IS : Set where
  ix : S → IS
  ar : S → IS

data Ctx : Set where
  ε : Ctx
  _>_ : Ctx → IS → Ctx

data _∈_ : IS → Ctx → Set where
  v_0 : is ∈ (Γ ▷ is)
  v_s : is ∈ Γ → is ∈ (Γ ▷ ip)

```

Note that while our contexts are non-dependent (*i.e.* the shapes do not depend on the terms), we use non-trivial shape dependencies within the constructors. The embedded language does not have a notion of shape as a value, therefore all the shape dependencies are handled by Agda, keeping our language simply typed (shaped). This separation is very helpful when it comes to writing embedded programs. We start with two helper definitions: a singleton shape that we call `unit` and the type for binary operations that we support (for now only addition and multiplication).

```

unit : S
unit = []

data Bop : Set where
  plus mul : Bop

```

The embedded language `E` includes: variables `var`; constants 0 and 1 given by `zero` and `one` correspondingly; three flavours of array constructor/eliminator pairs given by `imaps/sels`, `imap/sel` and `imapb/selb`; summation `sum`; conditional `zero-but` where the predicate is fixed to equality of two indices and the else branch is zero; `slide` and `backslide` exactly as described before; and numerical operations. The latter includes `logistic`, plus and multiplication, division by a constant `scaledown`, and unary `minus`. The definition of the embedded language `E` follows. We also introduce the syntax for infix plus and multiplication denoted `⊕` and `⊗` correspondingly.

```

data E : Ctx → IS → Set where
  var      : is ∈ Γ → E Γ is
  zero     : E Γ (ar s)
  one      : E Γ (ar s)

  imaps    : E (Γ ▷ ix s) (ar unit) → E Γ (ar s)
  sels     : E Γ (ar s) → E Γ (ix s) → E Γ (ar unit)

  imap     : E (Γ ▷ ix s) (ar p) → E Γ (ar (s ⊗ p))
  sel      : E Γ (ar (s ⊗ p)) → E Γ (ix s) → E Γ (ar p)

  imapb    : s * p ≈ q → E (Γ ▷ ix s) (ar p) → E Γ (ar q)
  selb     : s * p ≈ q → E Γ (ar q) → E Γ (ix s) → E Γ (ar p)

  sum      : E (Γ ▷ ix s) (ar p) → E Γ (ar p)
  zero-but : E Γ (ix s) → E Γ (ix s) → E Γ (ar p) → E Γ (ar p)

  slide    : E Γ (ix s) → s + p ≈ r → E Γ (ar r) → suc p ≈ u → E Γ (ar u)
  backslide : E Γ (ix s) → E Γ (ar u) → suc p ≈ u → s + p ≈ r → E Γ (ar r)

  logistic : E Γ (ar s) → E Γ (ar s)
  bin      : Bop → E Γ (ar s) → E Γ (ar s) → E Γ (ar s)
  scaledown : ℕ → E Γ (ar s) → E Γ (ar s)
  minus    : E Γ (ar s) → E Γ (ar s)
  let'     : E Γ (ar s) → E (Γ ▷ ar s) (ar p) → E Γ (ar p)

pattern _⊗_ a b = bin mul a b
pattern _⊕_ a b = bin plus a b

```

## 5.1 Reals

Explain that this is our module for abstract reals with operations and their properties and that we parametrise our evaluator with this module.

```

record Real : Set1 where
  field
    R : Set
    fromN : ℕ → R
    _+_ _*_ _÷_ : R → R → R
    _-e^_ : R → R
    logisticr : R → R
    logisticr x = fromN 1 ÷ (fromN 1 + e^ (- x))

```

## 5.2 Evaluation

The text in this section needs adjustment

We define the interpretation  $\llbracket \_ \rrbracket$  for  $(E \Gamma \text{ is})$  into the value  $(Val \text{ is})$  in the environment  $(Env \Gamma)$ . The values are either arrays or positions of the corresponding shape. Environments for the given context  $\Gamma$  are tuples of values of the corresponding shapes. The `lookup` function translates variables within the context into variables within the environment.

$Val : IS \rightarrow Set$	$Env : Ctx \rightarrow Set$	$lookup : is \in \Gamma \rightarrow Env \Gamma \rightarrow Val \text{ is}$
$Val (\text{ar } s) = Ar \ s \ R$	$Env \epsilon = \top$	$lookup \ v_0 \ (\rho, x) = x$
$Val (\text{ix } s) = P \ s$	$Env (\Gamma \triangleright is) = Env \Gamma \times Val \ is$	$lookup \ (v_s \ v) \ (\rho, \_) = lookup \ v \ \rho$

In the definition of  $\llbracket \_ \rrbracket$  we wrap the environment argument into double braces. This is an Agda-specific syntax for instance arguments<sup>4</sup> which behave similarly to hidden arguments, but they have a more powerful resolution algorithm. As a result we can omit mentioning the environment in recursive calls when it is passed unchanged.

```

llbracket _ \rrbracket : E \Gamma \text{ is} \rightarrow Env \Gamma \rightarrow Val \text{ is}
llbracket var \ x \rrbracket \rho = lookup \ x \ \rho
llbracket zero \rrbracket \rho = K (fromN 0)
llbracket one \rrbracket \rho = K (fromN 1)
llbracket imaps \ e \rrbracket \rho = \lambda i \rightarrow llbracket e \rrbracket (\rho, i) []
llbracket sels \ e \ e_1 \rrbracket \rho = K (llbracket e \rrbracket \rho (llbracket e_1 \rrbracket \rho))
llbracket imap \ e \rrbracket \rho = unnest \ \lambda i \rightarrow llbracket e \rrbracket (\rho, i)
llbracket sel \ e \ e_1 \rrbracket \rho = nest (llbracket e \rrbracket \rho) (llbracket e_1 \rrbracket \rho)
llbracket imapb \ m \ e \rrbracket \rho = Ar.imapb (\lambda i \rightarrow llbracket e \rrbracket (\rho, i)) m
llbracket selb \ m \ e \ e_1 \rrbracket \rho = Ar.selb (llbracket e \rrbracket \rho) m (llbracket e_1 \rrbracket \rho)
llbracket sum \ e \rrbracket \rho = Ar.sum (Ar.zipWith _+_ ) (K (fromN 0)) (\lambda i \rightarrow llbracket e \rrbracket (\rho, i))
llbracket zero-but \ i \ j \ e \rrbracket \rho = if llbracket i \rrbracket \rho  $\stackrel{?}{\leq}_p$  llbracket j \rrbracket \rho then llbracket e \rrbracket \rho else K (fromN 0)
llbracket slide \ e \ p \ e_1 \ s \rrbracket \rho = Ar.slide (llbracket e \rrbracket \rho) p (llbracket e_1 \rrbracket \rho) s
llbracket backslide \ e \ e_1 \ s \ p \rrbracket \rho = Ar.backslide (llbracket e \rrbracket \rho) (llbracket e_1 \rrbracket \rho) s (fromN 0) p
llbracket logistic \ e \rrbracket \rho = Ar.map logisticr (llbracket e \rrbracket \rho)

```

<sup>4</sup>For more details on instance arguments see: <https://agda.readthedocs.io/en/v2.6.3/language/instance-arguments.html>

```

491  $\llbracket e \boxplus e_1 \rrbracket \rho = \text{Ar.zipWith } \_+ \_ (\llbracket e \rrbracket \rho) (\llbracket e_1 \rrbracket \rho)$ 
492  $\llbracket e \boxtimes e_1 \rrbracket \rho = \text{Ar.zipWith } \_ * \_ (\llbracket e \rrbracket \rho) (\llbracket e_1 \rrbracket \rho)$ 
493  $\llbracket \text{scaledown } n \ e \rrbracket \rho = \text{Ar.map } (\_ \div \text{fromN } n) (\llbracket e \rrbracket \rho)$ 
494  $\llbracket \text{minus } e \rrbracket \rho = \text{Ar.map } \_ - \_ (\llbracket e \rrbracket \rho)$ 
495  $\llbracket \text{let } e_1 \rrbracket \rho = \llbracket e_1 \rrbracket (\rho, \llbracket e \rrbracket \rho)$ 

```

With the above definition we can better explain the choices of language constructors. The most important question to clarify is why do we have three array constructors/eliminators. As the only conceptual datatype of our language is an array (of some shape), we do not have any direct way to talk about array elements. Therefore, we model the type of array elements (scalars) as arrays of a singleton shape. As can be seen, scalar selection `sel`, returns a singleton array (application of `K`) where all the element(s) are equal to the element we are selecting. The corresponding array constructor `imap`, makes sure that if we compute  $s$  elements of the shape `unit`, we produce an array of shape  $s$  (and not  $s \otimes \text{unit}$ ). This solves the problem of constructing arrays from scalars, but how do we construct an array of a product shape? Given that we have an expression in the context  $(\Gamma \triangleright \text{ix } s \triangleright \text{ix } p)$ , we need to produce an array of  $s \otimes p$ . There are several ways how to solve this (e.g. introducing nest/unnest or projections and pairing on indices), but it is clear that we need something more than just an `imap`. This is the reason to introduce `imap/sel` pair which operates on arrays of product shapes. As average pooling operates on blocked arrays, we need a construction to express this in `E`. One could introduce explicit `block/unblock`, but we merge blocking/unlocking action with `imap/sel` obtaining `imapb/selb`. Our `sum` constructor gets an argument in the extended context which is summation index, so conceptually we generate the values at every summation index before summing these values together. As a result, we only need one instance of `sum` which makes our expressions a little tidier.

### 5.3 Weakening and Substitution

#### Adjust the text

As our language has explicit de Bruin variables (as opposed to HOAS [17] approaches), we need the means to do weakening and substitution when we optimise expressions in `E`. Our language is intrinsically typed(shaped) which makes the definition of both operations challenging. However, this problem has been well-understood, and we adopt the solution from [11]. We only show the basic mechanisms of the definition, for full details refer to [11].

The key structure that gives rise to weakening and substitution is a function that computes the context  $\Gamma$  *without* the variable  $v$  (denoted  $\Gamma / v$ ). Then we define the weakening for variables (`wkv`) and expressions (`wk`) that take a variable or expression in the context without the variable  $v$  and return this variable or expression in the context where  $v$  is present.

```

528 data  $\_ \subseteq \_ : \text{Ctx} \rightarrow \text{Ctx} \rightarrow \text{Set where}$ 
529    $\epsilon : \epsilon \subseteq \epsilon$   $\text{wkv} : \Gamma \subseteq \Delta \rightarrow is \in \Gamma \rightarrow is \in \Delta$ 
530    $\text{skip} : \Gamma \subseteq \Delta \rightarrow \Gamma \subseteq (\Delta \triangleright is)$   $\text{wk} : \Gamma \subseteq \Delta \rightarrow E \ \Gamma \ is \rightarrow E \ \Delta \ is$ 
531    $\text{keep} : \Gamma \subseteq \Delta \rightarrow (\Gamma \triangleright is) \subseteq (\Delta \triangleright is)$ 

```

#### Expand this (or hide?)

```

534  $\subseteq\text{-eq} : \Gamma \subseteq \Gamma$ 
535  $\subseteq\text{-eq } \{\epsilon\} = \epsilon$ 
536  $\subseteq\text{-eq } \{\Gamma \triangleright x\} = \text{keep } \subseteq\text{-eq}$ 

```

540  $\_↑ : E \Gamma \text{ is} \rightarrow E (\Gamma \triangleright ip) \text{ is}$

541  $\_↑ = \text{wk} (\text{skip} \subseteq \text{eq})$

542 Say something about substitution

544  $\text{data Sub } (\Gamma : \text{Ctx}) : \text{Ctx} \rightarrow \text{Set} \text{ where}$

545  $\epsilon : \text{Sub } \Gamma \epsilon$

546  $\_ \triangleright \_ : \text{Sub } \Gamma \Delta \rightarrow E \Gamma \text{ is} \rightarrow \text{Sub } \Gamma (\Delta \triangleright \text{is})$

548  $\text{sub} : E \Delta \text{ is} \rightarrow \text{Sub } \Gamma \Delta \rightarrow E \Gamma \text{ is}$

549 We can define identity substitution as follows:

551  $\text{sub-id} : \text{Sub } \Gamma \Gamma$

552  $\text{sub-id } \{\epsilon\} = \epsilon$

553  $\text{sub-id } \{\Gamma \triangleright x\} = \text{skeep sub-id}$

555 As our context do not encode explicit dependencies between the variables, we can easily define  
556 a substitution that swaps two top variables in the context. This will be used later for optimising  
557 programs in our DSL.

558  $\text{sub-swap} : \text{Sub } (\Gamma \triangleright \text{is} \triangleright ip) (\Gamma \triangleright ip \triangleright \text{is})$

559  $\text{sub-swap} = (\text{sdrop } (\text{sdrop sub-id}) \triangleright \text{var } v_0) \triangleright \text{var } (v_s \ v_0)$

561 *Syntax.*

562 Explain that we want to simplify the life of programers by introducing HOAS-like syntax,  
563 which is difficult, as our DSL is intrinsically-typed.

566  $\text{data Prefix} : (\Gamma \Delta : \text{Ctx}) \rightarrow \text{Set} \text{ where}$

567 *instance*

568  $\text{zero} : \text{Prefix } \Gamma \Gamma$

569  $\text{suc} : \{\{ \text{Prefix } \Gamma \Delta \} \} \rightarrow \text{Prefix } \Gamma (\Delta \triangleright \text{is})$

571 – A term that can be lifted into larger contexts

572  $\text{GE} : \text{Ctx} \rightarrow \text{IS} \rightarrow \text{Set}$

573  $\text{GE } \Gamma \text{ is} = \forall \{\Delta\} \rightarrow \{\{ \text{Prefix } \Gamma \Delta \} \} \rightarrow E \Delta \text{ is}$

574 – A variable that can be lifted into larger contexts

575  $\text{GVar} : \text{Ctx} \rightarrow \text{IS} \rightarrow \text{Set}$

576  $\text{GVar } \Gamma \text{ is} = \forall \{\Delta\} \rightarrow \{\{ p : \text{Prefix } \Gamma \Delta \} \} \rightarrow \text{is} \in \Delta$

577 – Lift var

578  $\text{V} : \text{is} \in \Gamma \rightarrow \text{GVar } \Gamma \text{ is}$

579  $\text{V } v \{\{ p = \text{zero} \} \} = v$

580  $\text{V } v \{\{ p = \text{suc} \} \} = v_s (\text{V } v)$

581 We can implement HOAS operators that will allow to use bound variables under further binders.

582 – Use GE GVar and V to define HOAS-style imap, imaps, and impab

583  $\text{Imap} : \forall \{\Gamma\}$

584  $\rightarrow (\text{GE } (\Gamma \triangleright \text{ix } s) (\text{ix } s) \rightarrow E (\Gamma \triangleright \text{ix } s) (\text{ar } p))$

```

589   → E Γ (ar (s ⊗ p))
590 Imap f = imap (f λ {Δ} { p } → var (V v0))
591
592 We do the same wrappers for sum, imaps, imapb, and the one for let', providing a familiar let-like
593 syntax.
594 Let-syntax : ∀ {Γ}
595   → (E Γ (ar s))
596   → (GE (Γ ▷ (ar s)) (ar s) → E (Γ ▷ (ar s)) (ar p))
597   → E Γ (ar p)
598 Let-syntax x f = let' x (f λ {Δ} { p } → var (V v0))
599
600 syntax Let-syntax e (λ x → e') = Let x := e In e'

```

The final convenience operator that we are missing is the ability to represent contexts in the HOAS style. First of all, we define a helper function that appends a list of **IS** at the end of some context  $\Gamma$ :

```

605 ext : Ctx → List IS → Ctx
606 ext Γ [] = Γ
607 ext Γ (x :: l) = ext (Γ ▷ x) l
608

```

Then, we define **lfun** that for the given list of **IS**-es ( $l = [is_1, \dots, is_n]$ ), some context  $\Gamma$  and some **IS**  $ip$  computes an Agda function of type  $(GE (ext \Gamma l) is_1 \rightarrow \dots \rightarrow GE (ext \Gamma l) is_n \rightarrow E (ext \Gamma l) ip)$ . The function **lvar** lifts a variable in some context  $\Gamma$  into an extended context.

```

612 lfun : (l : List IS) (Γ : Ctx) (is : IS) → Set
613 lvar : ∀ l → is ∈ Γ → GE (ext Γ l) is

```

With these helper functions we define the **Lcon** helper that for the given list of types  $l$ , resulting type  $is$ , the initial context  $\Gamma$  and the function of type **lfun**  $l \Gamma is$  computes the expression in the context **ext**  $\Gamma l$ .

```

618 Lcon : ∀ l is Γ → (f : lfun l Γ is) → E (ext Γ l) is
619 Lcon [] is Γ f = f
620 Lcon (x :: l) is Γ f = Lcon l is (Γ ▷ x) (f (lvar l v0))
621

```

This means that we can bind the last  $n$  elements of the context to Agda variables and use them safely under binders. For example, consider this expression:

```

624 _ : E _
625 _ = Lcon (ar (ι 5) :: ar (5 :: 5 :: []) :: []) (ar []) ε
626   λ a b → Sum λ i → sels a i sels (sel b i) i
627

```

where we  $a$  and  $b$  are bound to the arguments of the Agda's lambda term and which are used when computing expression in the context  $(\epsilon \triangleright 5 :: [] \triangleright 5 :: 5 :: [])$ .

*Primitives.* We are defining primitives that are needed for expresison our running example in  $E$ . We consider an example for the **conv**-olution:

```

632 conv : ∀ {Γ} → E Γ (ar r) → { s + p ≈ r } → E Γ (ar s) → { suc p ≈ u }
633   → E Γ (ar u)
634 conv f { s+p } g { ss }
635   = Sum λ i → (slide i s+p ⟨ f ⟩ ss) ⋈ Imaps λ j → sels ⟨ g ⟩ i
636
637

```

```

638 mconv : { s + p ≈ r } → (inp : E Γ (ar r)) (ws : E Γ (ar (u ⊗ s)))
639       (b_v : E Γ (ar u)) → { suc p ≈ w } → E Γ (ar (u ⊗ w))
640 mconv { sp } inp w_v b_v { su } =
641   Imap λ i → conv ⟨ inp ⟩ (sel ⟨ w_v ⟩ i) ⊞ Imaps λ _ → sels ⟨ b_v ⟩ i
642
643 avgp2 : ∀ m n → (a : E Γ (ar (m ℕ.* 2 :: n ℕ.* 2 :: [])))
644       → E Γ (ar (m :: n :: []))
645 avgp2 m n a =
646   Imaps λ i → scaledown 4 $ Sum λ j → sels (selb it ⟨ a ⟩ i) j
647
648 sqerr : (r o : E Γ (ar [])) → E Γ (ar [])
649 sqerr r o = scaledown 2 ((r ⊞ (minus o)) ⊞ (r ⊞ (minus o)))
650
651 meansqerr : (r o : E Γ (ar s)) → E Γ (ar [])
652 meansqerr r o = Sum λ i → sqerr (sels ⟨ r ⟩ i) (sels ⟨ o ⟩ i)
653
654 where ⟨ ⟩ lifts an expression into a generalised expression simply by applying weakening according
655 to the prefix.
656 Finally, the CNN embedded in E is given as follows:
657
658 cnn : E _ _
659 cnn = Lcon ( ar (28 :: 28 :: []) :: ar (6 :: 5 :: 5 :: [])
660           :: ar (6 :: []) :: ar (12 :: 6 :: 5 :: 5 :: [])
661           :: ar (12 :: []) :: ar (10 :: 12 :: 1 :: 4 :: 4 :: [])
662           :: ar (10 :: []) :: ar (10 :: 1 :: 1 :: 1 :: 1 :: [])
663           :: []
664           (ar []) ε
665 λ inp k1 b1 k2 b2 fc b target →
666 Let c11 := mconv inp k1 b1 In
667 Let c1 := logistic c11 In
668 Let s1 := (Imap {s = 6 :: []} λ i → avgp2 12 12 (sel c1 i)) In
669 Let c21 := mconv s1 k2 b2 In
670 Let c2 := logistic c21 In
671 Let s2 := (Imap {s = 12 :: 1 :: []} λ i → avgp2 4 4 (sel c2 i)) In
672 Let o1 := mconv s2 fc b In
673 Let o := logistic o1 In
674 Let e := meansqerr target o In
675 e
676
677
678
679
680
681
682
683
684
685
686

```

The code in the next section has not been adjusted yet, so do not change the 'lagda' files, it won't compile!

## 6 Automatic Differentiation

We implement automatic differentiation in reverse mode for expressions in E. We focus on reverse mode because it is of most interest in machine learning, and it is more challenging to implement.

We start with a brief introduction of the AD, for much more in-depth explanations refer to [3, 4]. Consider differentiating a function composition consisting of three functions:

$$y = (f \circ g \circ h) x$$

rewrite it using temporary variables:

$$\begin{aligned} w_0 &= x \\ w_1 &= h w_0 \\ w_2 &= g w_1 \\ w_3 &= f w_2 = y \end{aligned}$$

The chain rule gives us  $\frac{\partial y}{\partial x} = \frac{\partial y}{\partial w_2} \frac{\partial w_2}{\partial w_1} \frac{\partial w_1}{\partial x}$ . The difference between the forward and reverse mode lies in the direction that we traverse the chain rule. In forward mode we traverse the chain inside-out, and the reverse mode traverses the chain outside-in thus computing recursive relation:  $\frac{\partial y}{\partial w_i} = \frac{\partial y}{\partial w_{i+1}} \frac{\partial w_{i+1}}{\partial w_i}$ . For our example, we compute  $\frac{\partial y}{\partial w_2}$ , then  $\frac{\partial w_2}{\partial w_1}$  and finally  $\frac{\partial w_1}{\partial x}$ . While there seem to be no difference for functions of one variable, there is a big difference for functions of  $n$  variables as we can compute derivatives of all the non-dependent variables simultaneously. Consider an example of the  $z = f x y = \sin(xy + x)$ :

$$\begin{aligned} w_0 &= x \\ w_1 &= y \\ w_2 &= w_1 w_2 \\ w_3 &= w_2 + w_0 \\ w_4 &= \sin w_3 = z \end{aligned}$$

We compute the adjoints  $\bar{w}_i = \frac{\partial y}{\partial w_i}$  using the following rule. If  $w_i$  has successors in the computational graph, we can apply the chain rule as follows:

$$\bar{w}_i = \sum_{j \in \text{succ } i} \bar{w}_j \frac{\partial w_j}{\partial w_i}$$

For our example:

$$\begin{aligned} \bar{w}_4 &= 1 = \frac{\partial z}{\partial z} \\ \bar{w}_3 &= \bar{w}_4 \cos w_3 \\ \bar{w}_2 &= \bar{w}_3 \cdot 1 \\ \bar{w}_1 &= \bar{w}_2 w_0 \\ \bar{w}_0 &= \bar{w}_3 + \bar{w}_2 w_1 \end{aligned}$$

If we inline all the  $\bar{w}_i$  definitions and inspect the values of partial derivatives with respect to  $x$  and  $y$  we obtain expected results:  $\frac{\partial z}{\partial x} = \cos(xy + x)(y + 1)$  and  $\frac{\partial z}{\partial y} = \cos(xy + x)x$ .

#### Old text

In the implementation of the AD for **E** in some context  $\Gamma$ , we would like to obtain all the partial derivatives with respect to the variables in context  $\Gamma$ . Each partial derivative is itself an expression **E** in context  $\Gamma$ . Therefore, we need to define a data type for an environment of  $\Gamma$ -many expressions in context  $\Gamma$ . We call this environment **Env** defined as follows:

```
data Env : Ctx → Ctx → Set where
  e  : Env → Γ
```



```

736 skip : Env  $\Gamma \Delta \rightarrow$  Env  $(\Gamma \triangleright \text{ix } s) \Delta$ 
737  $\_ \triangleright \_ :$  Env  $\Gamma \Delta \rightarrow$  E  $\Delta (\text{ar } s) \rightarrow$  Env  $(\Gamma \triangleright \text{ar } s) \Delta$ 
738
739 data EE : Ctx  $\rightarrow$  Ctx  $\rightarrow$  Set where
740   env : Env  $\Gamma \Delta \rightarrow$  EE  $\Gamma \Delta$ 
741   let' : E  $\Delta (\text{ar } s) \rightarrow$  EE  $\Gamma (\Delta \triangleright \text{ar } s) \rightarrow$  EE  $\Gamma \Delta$ 

```

Note that **Env** only keeps array expressions, as (i) derivatives for indices do not exist; and (ii) we can always make an initial environment by populating all the elements with **zeros**.

Old text

We define several helper operations to manipulate environments: **env-zero** is an environment where all the values are **zeros**; **update** modifies the expression at the  $v$ -th position by applying  $f$  to it; **env-map** applies the function  $f$  from **E** to **E** to all the elements of the environment; and **env-zipWith** applies the binary function  $f$  on two environments point-wise. The types of these helper functions follow. As environments are very similar to lists, the implementation of the above functions are straight-forward.

```

752 ee-wk      :  $\Delta \subseteq \Phi \rightarrow$  EE  $\Gamma \Delta \rightarrow$  EE  $\Gamma \Phi$ 
753 ee-wk-zero : EE  $\Gamma \Delta \rightarrow$   $\Gamma \subseteq \Phi \rightarrow$  EE  $\Phi \Delta$ 
754 ee-tail    : EE  $(\Gamma \triangleright \text{is}) \Delta \rightarrow$  EE  $\Gamma \Delta$ 
755 zero-ee    : EE  $\Gamma \Delta$ 
756 ee-plus    :  $(\rho \rho : \text{EE } \Gamma \Delta) \rightarrow$  EE  $\Gamma \Delta$ 
757 ee-map-sum : EE  $\Gamma (\Delta \triangleright \text{ix } s) \rightarrow$  EE  $\Gamma \Delta$ 
758 ee-update+ : EE  $\Gamma \Delta \rightarrow$   $(v : \text{ar } s \in \Gamma) (t : \text{E } \Delta (\text{ar } s)) \rightarrow$  EE  $\Gamma \Delta$ 
759  $\_ \triangleright 0 :$  EE  $\Gamma \Delta \rightarrow$  EE  $(\Gamma \triangleright \text{ar } s) (\Delta \triangleright \text{ar } s)$ 
760

```

We define the function  $\nabla$  that takes an expression **E** and the seed which is the multiplier on the left of the chain, and we compute a function from that updates the environment.

```

763 {-# TERMINATING #-}
764  $\nabla_l : \text{E } \Gamma (\text{ar } s) \rightarrow$  EE  $(\Gamma \triangleright \text{ar } s) \Gamma \rightarrow$  EE  $\Gamma \Gamma$ 
765  $\nabla \Sigma : (e s : \text{E } (\Gamma \triangleright \text{ix } s) (\text{ar } p)) \rightarrow$  EE  $\Gamma \Gamma \rightarrow$  EE  $\Gamma \Gamma$ 
766
767  $\nabla : (e s : \text{E } \Gamma \text{ is}) \rightarrow$  EE  $\Gamma \Gamma \rightarrow$  EE  $\Gamma \Gamma$ 
768  $\nabla \{ \text{is} = \text{ix } \_ \} (\text{var } x) s = \text{id}$ 
769  $\nabla \{ \text{is} = \text{ar } \_ \} (\text{var } x) s = \lambda \delta \rightarrow \text{ee-update+ } \delta x s$ 
770  $\nabla \text{ zero } s = \text{id}$ 
771  $\nabla \text{ one } s = \text{id}$ 
772
773  $\nabla (\text{imaps } e) s = \nabla \Sigma e (\text{sels } (s \uparrow) (\text{var } v_0))$ 
774  $\nabla (\text{imap } e) s = \nabla \Sigma e (\text{sel } (s \uparrow) (\text{var } v_0))$ 
775  $\nabla (\text{E.imapb } m e) s = \nabla \Sigma e (\text{E.selb } m (s \uparrow) (\text{var } v_0))$ 
776
777  $\nabla (\text{sels } e i) s = \nabla e (\text{imaps } (\text{zero-but } (\text{var } v_0) (i \uparrow) (s \uparrow)))$ 
778  $\nabla (\text{sel } e i) s = \nabla e (\text{imap } (\text{zero-but } (\text{var } v_0) (i \uparrow) (s \uparrow)))$ 
779  $\nabla (\text{E.selb } m e i) s = \nabla e (\text{E.imapb } m (\text{zero-but } (\text{var } v_0) (i \uparrow) (s \uparrow)))$ 
780
781  $\nabla (\text{E.sum } e) s = \nabla \Sigma e (s \uparrow)$ 
782  $\nabla (\text{zero-but } i j e) s = \nabla e (\text{zero-but } i j s)$ 
783  $\nabla (\text{E.slide } i p e su) s = \nabla e (\text{E.backslide } i s su p)$ 
784

```

```

785  $\nabla (\text{E.backslide } i \ e \ su \ p) \ s = \nabla \ e \ (\text{E.slide } i \ p \ s \ su)$ 
786
787  $\nabla (e \boxplus e_1) \quad s = \nabla \ e \ s \circ \nabla \ e_1 \ s$ 
788  $\nabla (e \boxtimes e_1) \quad s = \nabla \ e \ (s \boxtimes e_1) \circ \nabla \ e_1 \ (s \boxtimes e)$ 
789  $\nabla (\text{scaledown } x \ e) \ s = \nabla \ e \ (\text{scaledown } x \ s)$ 
790  $\nabla (\text{minus } e) \quad s = \nabla \ e \ (\text{minus } s)$ 
791  $\nabla (\text{logistic } e) \quad s = \nabla \ e \ (\text{let}' \ (\text{logistic } e) \ ((s \uparrow) \boxtimes \text{var } v_0 \boxtimes (\text{one} \boxplus \text{minus } (\text{var } v_0))))$ 
792
793  $\nabla (\text{let}' \ e \ e_1) \quad s = \lambda \ \delta \rightarrow \nabla_l \ e \ (\text{let}' \ e \ (\nabla \ e_1 \ (s \uparrow) \ (\delta \triangleright 0)))$ 
794
795  $\nabla \Sigma \ e \ s \ \delta = \text{ee-plus } \delta \ \$ \text{ ee-tail } \$ \text{ ee-map-sum } (\nabla \ e \ s \ \text{zero-ee})$ 
796
797  $\nabla_l \ e \ (\text{env } (\rho \triangleright x)) = \text{ee-tail } \$ \text{ let}' \ x \ (\nabla \ (e \uparrow) \ (\text{var } v_0) \ (\text{env } \rho \triangleright 0))$ 
798  $\nabla_l \ e \ (\text{let}' \ x \ \rho) = \text{let}' \ x \ (\text{ee-tail } \$ \nabla_l \ (e \uparrow) \ (\text{ee-wk-zero } \rho \ (\text{keep } (\text{skip } \sqsubseteq \text{-eq}))))$ 
799

```

### Old text

Let us now walk through the cases. Derivative of constants (**zero** and **one**) is zero, so nothing needs to be updated in the environment. Index variables are not stored in the environment, so no updates are needed either. If we differentiate the variable  $x$  with some seed  $s$ , we update the  $x$ -th position in the environment by adding  $s$  to it. Differentiation of **imaps** proceeds as follows: we recursively apply  $\nabla$  to  $e$  (in the context  $\Gamma \triangleright (\text{ix } p)$ ) with the element of the original seed  $s$  selected at the top variable. This gives us the environment in the extended context, then we map **sum** to every element of the environment to accumulate the derivatives at every index. When differentiating selections we recurse on the array we are selecting from with the seed that is zero everywhere except the index we were selecting at. Differentiating conditional is straight-forward, as  $i$  and  $j$  must be in the context, we can simply differentiate  $e$  with the condition on seed. If indices were equal, we will compute the update, otherwise we will differentiate with seed **zero** which has no effect. As we are operating in a total language, there is no need to worry about pulling the expression out of conditional. The argument of **sum** lives in the extended context, so we apply the same rules as for the **imap** family, except we propagate the original seed to all the summands. Addition and multiplication rules are straight-forward application of differentiation rules. The **slide/backslide** pair forms a satisfying  $\nabla$ -symmetry. Finally, **scaledown**, **minus** and **logistic** follow the rules of differentiation.

## 6.1 Optimisation

Replace the old text, explain that we are semantically-preserving now

Our algorithm often delivers expressions that are not computationally efficient. While we can hope for the backend to take care of this, it is relatively easy to implement a number of rewriting rules prior to extraction. We constructed **E** such that no computation is happening in the shape or context positions. As a result, dependent pattern-matching is always applicable on **E** expressions, and our optimisations can be formulated very concisely. We omit constant-folding like rewrites such as addition with zero and multiplication by one and focus on less trivial cases that have to do with selections and sum. Consider the snippet of the optimiser for **sel<sub>i</sub>** and **sum**.

## 6.2 Compilation

We had two reasons to define the embedded language **E**. Firstly, **E** makes it possible to implement automatic differentiation within Agda, as we described in the previous section. Secondly, we compile

expressions in **E** into a programming language that can produce efficient code. This section describes extraction process into Futhark.

Futhark is a functional language with automatic memory management and built-in type for arrays. Futhark provides key array combinators such as `map` and `reduce`, which makes translation process straight-forward. There are two non-trivial aspects of the process that we describe below.

*Static Ranks.* In Futhark, array ranks are static. This means that it is not possible to translate any expression in **E** into Futhark. We assume that all the ranks are known statically, which is true for many numerical applications including our running example.

*Normalisation.* Consider translating an expression like `sel (imap  $\lambda i \rightarrow e$ ) u`. If you were to treat arrays as functions and selections as applications, then the above expression can be normalised into `e[i := u]`. One could hope that Futhark could do such a  $\beta$ -reduction on the generated code, but this is not the case. The intuition for this choice is that in Futhark arrays are tabulated functions, and inlining arbitrary evaluation of array elements may have a significant performance cost. For example, in the expression `let a = imap  $\lambda i \rightarrow e$  in imap  $\lambda j \rightarrow a[f j]$` , Futhark allocates memory for `a` and computes all the values, and within the body of the `let`, selection actually looks up the elements. If we were to inline `a` by replacing `a[f j]` with `e[i := f j]`, we loose sharing by potentially recomputing `e` much more often than needed (e.g. assume that `i` ranges over 10 elements, but `j` over  $10^5$ ). Resolving when such inlining is beneficial for performance is non-trivial, therefore Futhark (and many other array languages) do not inline computation of array elements. For our running example, naive translation results in too many cases when arrays are constructed just to select an element from them. Therefore, we need some notion of normalisation prior extraction.

We are going to combine normalisation and extraction in a single step, resulting in something similar to normalisation by evaluation. We model Futhark arrays as Agda functions space which makes it easy to encode normalisation steps. Futhark indices for an array of shape `s` are given by the type `Ix` which is simply a list of strings (the name of the index) per dimension:

```
data Ix : S → Set where
  [] : Ix []
  _::_ : String → Ix s → Ix (n :: s)
```

The `Sem` function gives an interpretation to types of **E** expressions. Indices are just interpreted as `Ix` of the corresponding shape. Array types are morally functions from indices to strings. However, in the definition the type is a little more complicated:

```
Sem : IS → Set
Sem (ar s) = (Ix s → State ℕ ((String → String) × String))
Sem (ix s) = Ix s
```

Let us explain the complexity of the array type. First of all, the codomain of the array is wrapped into a state monad which gives a source of fresh variable names. Within the monad we have a pair we have a function which represents a context for the actual array computation which is the second argument. This context is needed because of the interplay between `let` bindings and `imap`. Consider for a moment that we do not have explicit context in the type for `ar` and we are compiling an expression `Let z := zero in Imap  $\lambda i \rightarrow z$`  which can result in something like:

```
f : Ix s → State ℕ String
f i = return ("let z = 0 in " ++ ( $\lambda j \rightarrow$  "z") i)
```

If we select into this array (by applying it to some index expression) or compose it with other functions, this works as expected. However, at a certain point we may need to turn this expression

into the actual array, which looks something like `"imap λ i → " ++ f "i"`. However, this expression evaluates to `"imap λ i → let z = 0 in z"`, but this inlines computation of let binding in the body of the `imap`, which may have a serious performance impact if let binds a non-trivial computation. By introducing contexts in `Sem`, we just control where the `imap` code is injected. Generally speaking, our strategy here is to preserve sharing that is introduced by let bindings, yet normalise bound expressions and bodies.

For the actual extraction we define the environment of Futhar values called `FEnv`. Two functions that actually do the translation are `to-fut` which computes the `Sem` value, and `to-str` that calls `to-fut` and wraps the result with `imap` as we described above.

```

FEnv : Ctx → Set
FEnv ε = ⊤
FEnv (Γ ▷ is) = FEnv Γ × Sem is

to-fut : E Γ is → FEnv Γ → State ℕ (Sem is)
to-str : E Γ (ar s) → FEnv Γ → State ℕ String

```

Consider two cases of `to-fut` for `imap` and `sel`. In both cases the array we are constructing or selecting from is of shape  $s + p$ . We use two helper functions `ix-curry` and `ix-uncurry` that translate between funtions of type  $\text{Ix } (s ++ p) \rightarrow X$  and  $\text{Ix } s \rightarrow \text{Ix } p \rightarrow X$ . In the `imap` case we generate a function keeping potential let chains within the `imap` expression. In case of `sel`, we are computing the value of the array we are selecting from (i.e.  $a$ ) and within the returned expression we apply  $a$  to the corresponding indices — this is normalisation step.

```

to-fut (imap {s = s} e) ρ =
  return $ ix-uncurry {s} λ i j → do
    b ← to-fut e (ρ, i)
    f, b' ← b j
    return (id, f b')

to-fut (sel e e₁) ρ = do
  a ← to-fut e ρ
  i ← to-fut e₁ ρ
  return λ j → do
    f, a' ← ix-curry a i j
    return (f, a')

```

The rest of the code generator looks very similar, therefore we omit it here but the full code is available in the supplementary materials.

The code in the next section has not been adjusted yet, so do not change the 'lagda' files, it won't compile!

## 7 Performance

One of the goals of this work is to demonstrate that it is possible to formulate the problem in a proof assistant and then pass it on to the other system that can run the algorithm efficiently. In order to substantiate this claim, we compare the running times of the code that we generate from the specification at the end of the Section 6 and the hand-written SaC code from [23]. We are not interested in an exhaustive performance study similar to what is provided in [23]. Instead, we take the version from that paper as reference point and we aim to find out whether we are in the same ballpark.

We take the code from [23], make sure that it runs, and we replace the hand-written CNN training with the Agda-generated one. Our first observation is that both versions produce the same results, and none of the shape constraints that we defined in Section ?? fired. This means that our code generation is working. Unfortunately, the runtime comparison revealed that our version is about 10× slower than the hand-written SaC version.

We got in touch with the SaC team who provided a lot of support in identifying the causes of the disappointing difference in performance. It turns out that the main culprit has to do with the inability to optimise away selections into tensor comprehensions in a few situations. With-Loop-Folding [19],

SaC's mechanism for fusing tensor comprehensions fails to fold tensor comprehensions that are nested and cannot be flattened statically. In simple terms, the expression  $\{iv \rightarrow e(iv)\}[jv]$  in some situation does not reduce to  $e(jv)$  which, in our generated code, is essential to match the hand-written performance. As a result, several arrays were created simply to make a single selection into them. The original code never ran into this problem as the hand-written code avoided such patterns. Our [E](#) optimiser from Section 6.1 could not help either, because the problem was occurring after the SaC primitives such as `slide` and `block` were inlined.

After numerous attempts on altering [E](#) to fit SaC requirements and the SaC team trying to implement some of the missing optimisations, on the February 23rd (6 days before the ICFP deadline) we realised that the best runtime we can obtain is still  $6\times$  slower than the hand-written code. The compiler is too sensitive to the flavour of the code that we pass to it, and when certain patterns are not recognised, there is very little one can do other than trying to fit those patterns. However, this is not always possible with the generated code. Performance is frustrating!

## 7.1 Generating C

After overcoming the natural instinct to give up, we realised that the real power of the proposed approach lies in the ability to modify any part of the code generation pipeline. This includes swapping the back-end language of choice to something else. Therefore, we decided to try generating C code instead of SaC code.

While C is a canonical low-level language, it has excellent support for multi-dimensional arrays, given that the ranks are statically known. At runtime these arrays are flattened vectors, they do not have to live on the stack, and the language takes care of multi-dimensional indexing.

However, the key difference between the C and SaC is memory management. SaC is a functional language that uses reference counting to automate operations on allocating and freeing memory. In C these decisions are manual, and as we have seen before, excessive memory allocation is detrimental for the runtime. For our use case we avoid memory management problem entirely, by assuming that all the variables in the [Chain](#) have to be preallocated, and if we need any temporary array variables when extracting array values, we fail extraction. This way we guarantee that no memory allocation is ever needed.

Meeting such a requirement means that we need to optimise away operations like `slide/backslide` as they require conceptual array allocation. The same goes for `imaps` appearing in some of the argument positions. Putting these considerations together, we extended [E](#) with the following explicit operations on indices:

```

data E : Ctx → IS → Set where
  div      : s * p ≈ q → (i : E Γ (ix q)) → E Γ (ix s)
  mod      : s * p ≈ q → (i : E Γ (ix q)) → E Γ (ix p)
  ix-plus   : (i : E Γ (ix s)) → (j : E Γ (ix u)) → suc p ≈ u → s + p ≈ r → E Γ (ix r)
  ix-minus  : (i : E Γ (ix r)) → (j : E Γ (ix s)) → s + p ≈ r → suc p ≈ u
              → (e : E (Γ ▷ ix u) (ar q)) → E Γ (ar q)
  ix-minusr : (i : E Γ (ix r)) → (j : E Γ (ix u)) → s + p ≈ r → suc p ≈ u
              → (e : E (Γ ▷ ix s) (ar q)) → E Γ (ar q)
  - ...

```

The `div` and `mod` constructors perform point-wise division or modulo operation on the index  $i$  and the shape  $p$ . This is needed to express selections into blocked arrays as we have seen in Section ?? . The `ix-plus` is a point-wise addition of  $i$  and  $j$ . The `ix-minus` and `ix-minusr` correspond to left and right subtraction from the Section 4. The meaning of these constructors is follows: if  $j$  can be

subtracted from  $i$  (in the sense of existence of inverse to  $\oplus_p$  exists) then we evaluate  $e$  at that index, otherwise we return zero.

**7.1.1 Optimisations.** We add the following optimisations to facilitate removal of temporary arrays in the generated code. We show the only ones that we added, all the optimisations we defined before are still valid.

```

opt : E  $\Gamma$  is  $\rightarrow$  E  $\Gamma$  is
opt (sels e e1) with opt e | opt e1
... | imapb m e | i = sels (sub v0 e (div m i)) (mod m i)
... | slide i pl a su | k = sels a (ix-plus i k su pl)
— | ... as before ...

```

Here we optimise away scalar selections into blocked imaps. Recall that  $m$  tells us that we have an array of shape  $s * p$ , and  $e$  computes blocks of shape  $p$ . If we are selecting into such a blocked array at the index  $i$ , we know that we are selecting  $(i/p)$ -th block, and from that block we are selecting  $(i \% p)$  element. Existence of explicit **div** and **mod** operations on indices makes it possible to implement this rewrite rule that is again very similar to  $\beta$ -reduction.

```

opt (sum e) with opt e
... | zero-but (var i) (ix-plus (var j) (var k) su pl) a with eq? v0 i | eq? v0 j | eq? v0 k
... | neq _ i' | neq _ j' | eq = ix-minus (var i') (var j') pl su a
... | neq _ i' | eq | neq _ k' = ix-minusr (var i') (var k') pl su a
... | _ | _ | _ = sum (zero-but (var i) (ix-plus (var j) (var k) su pl) a)
— | ... as before ...

```

Here we are dealing with the sum over summation index  $t$  where the inner expression is a conditional on indices of the form  $i == j + k ? e : \emptyset$ . Here we apply the same comparison of index variables as before. If  $k$  happens to be the variable  $t$ , then overall sum will only add one non-zero element at  $(i - j)$ -th index, given that this (left) subtraction is possible in the sense of existence of the inverse to  $\oplus_p$  operation defined in Section 4.2. The same happens when the summation index  $t$  is equal to  $j$ , we only need to consider  $(i - k)$ -th element given that this (right) subtraction is possible. One could cover other cases where  $t$  is equal to  $i$ , or when  $i$  and  $j + k$  are swapped, but these are not occurring in our running example.

```

opt (scaledown x e) with opt e
... | sum a = sum (scaledown x a)
— | ... as before ...

```

Finally, here is a rule that looks very innocent in the high-level language, yet becomes of importance in the low-level one. The rule says that if we are summing the array and then dividing it by a constant, we should move division inside the summation. The reason for this rewrite rule being important is when the result of the sum is non-scalar, we need to create a temporary array, before scaling down all its elements. A language with first class arrays can obviously take care of such minor details, but in C we have to be explicit about it.

**7.1.2 Code Generation.** Due to space limitations, we only consider the basic mechanisms we used in the code generator, all the code is available in supplementary materials. We use heap-allocated multi-dimensional arrays that can be defined as follows:

```
float(*k1)[6][5][5] = malloc(sizeof(*k1));
```

This ensures that  $k1$  is represented as a continuous region of memory of size  $6 * 5 * 5$  floats. When such arrays are indexed (e.g.  $(*k1)[i][j][k]$ ), the indices are translated into a single offset into the continuous memory. Therefore, there is no pointer chasing which makes this approach efficient at runtime. As C uses row-major order to compute the offsets, we do obtain partial array selections on the left, e.g.  $(*k1)[i]$  is a  $5 \times 5$  array that can be further indexed or passed to `sizeof` that correctly identifies the size of this subarray. Surely, this is a pointer into the  $k1$  array, so all the modifications to  $(*k1)[i]$  will modify  $k1$ . As a great convenience feature, C compiler tracks the ranges of the indices and produces warnings in cases when it figures out that ranges of indices and the array we are indexing do not match.

Whenever we translate some  $e$  in **E** into C, we have to provide a storage where  $e$  has to be written to. In case of compiling the **Chain** every local variable becomes such a storage for the bound expression. Therefore, our extractor always has a result variable as an argument.

For example, let us consider an expression  $a \boxplus a$  of shape  $(l\ 5 \otimes l\ 5)$ , where  $a$  is mapped to the C variable `float (*a)[5][5]` that is written to the result variable `float (*r)[5][5]`. Here is the code that we generate:

```
for (size_t x1_1 = 0; x1_1 < 5; x1_1++) {
  for (size_t x1_2 = 0; x1_2 < 5; x1_2++) {
    (*r)[x1_1][x1_2] = ((*a)[x1_1][x1_2] + (*a)[x1_1][x1_2]);
  }
}
```

We started with checking that  $a \boxplus a$  is a *selectable* expression. This means that we can always generate expression at the given index. As we know that the shape of  $a \boxplus a$  is  $(l\ 5 \otimes l\ 5)$ , we need to generate a loop nest of that shape that assigns where we assign the expression at the given index to the result at the given index.

We need to distinguish whether we are writing into the result or adding into it as in cases when dealing with **sum**. Consider the code that is generated for `(sum (sels (a (var  $v_0$ ))))` where we are adding all the elements of the array  $a$  into result variable `float (*r)[1]`:

```
for (size_t x2_1 = 0; x2_1 < 5; x2_1++) {
  for (size_t x2_2 = 0; x2_2 < 5; x2_2++) {
    for (size_t x3_1 = 0; x3_1 < 1; x3_1++) {
      (*r)[x3_1] += (&(*a)[x2_1][x2_2])[x3_1];
    }
  }
}
```

Two things are happening here, first we generate `+=` assignment and we make an implicit assumption that resulting variables are initialised to zero. In the extractor, additionally to the resulting variable we track whether we need to do an assignment or assignment with addition. Secondly, while  $a$  is two-dimensional, we have three-dimensional loop nest. The latter comes from the representation of scalars as 1-element vectors. When we resolved the two-dimensional summation index  $x2$ , we know that we need to assign into the object of shape  $(l\ 1)$ , but the left-hand-side is a scalar (float). The trick here is that in C we can always turn scalars into 1-element vectors by simply taking the address of the scalar. This is why we have this 1-iteration for-loop over  $x3_1$  that will be immediately optimised away by the C compiler.

Finally, when it comes to the operation on indices, such as addition, subtraction, division or modulo, we generate the corresponding operation on the individual loop indices.

Remaining details of the code generation take care of traversing through the structure of **E** with some plumbing that has to do with generating loop-nests around expressions and checking that they are selectable.

**7.1.3 Running the Generated C Code.** In order to run the generated C code we translate the boilerplate code from SaC to C. While doing so, we made sure that our code can be run in parallel. While the SaC compiler does this automatically, there is one obvious loop that requires parallelisation



which is computation of the batch. When we train the CNN, we take a batch of images and the weights and we compute gradients for those weights per every image. After that we average all the gradients in the batch, and we update the weights, after all the batch is processed. Clearly, all the gradient computations in the batch can run in parallel. We achieve this by organising the batch loop such that all the gradients are stored in a separate memory region, and we parallelise this loop using OpenMP annotations.

We verify that the code that we generate compute the same results as the hand-written SaC code. Then we replicate the experiment from the [23] using 40 epochs, 100 images in the batch, and feeding 10000 training images. We run the experiment on the 18-core 13th Gen Intel(R) Core(TM) i5-13600K machine using sac2c version 1.3.3-MijasCosta-1161-gb543c and the GCC compiler version 12.2.0. The first thing that we learn is that our generated C code is sensible (factor of 3 running time) to the compilation flags that we enable. We identified the set of flags that when passed to both compilers<sup>5</sup>, the runtime at the largest number of cores are 11s for the hand-written SaC implementation and 13.5s for the generated C code, with very little variance. This 20% difference is orthogonal to parallel execution, as it is also observed when running the code on a single core. The set of flags has to do with floating point operations: `-fno-signed-zeros` ignores the distinction between negative and positive zeroes that is given by IEEE 754 standard, allowing to reduce  $(-0.0 \times x)$  to 0.0; `-fno-math-errno` does not set `errno` after calling math functions; `-fno-trapping-math` and `-fassociative-math` make sure that we can assume associativity of floating point operations which does not hold according to the IEEE 754.

The main performance difference comes from the fact that compiled SaC code uses less intermediate arrays, significantly reducing the number of memory writes. There are numerous ways how to improve the performance of the generated C code, but for the purposes of this paper we consider that getting within 20% of the hand-written SaC code is sufficient evidence for our hypothesis that the two-languages approach seems viable for achieving proved correctness and performance. We have automatic differentiation in the safe environment that generates the C code that runs almost as fast as the hand-written SaC code.

## 8 Related Work

## 9 Conclusions

The paper demonstrates a technique of developing high performance applications with strong correctness guarantees. The key insight lies in using a proof assistant in cooperation with a high-performance language of choice. This gives a clear separation of concerns that is very difficult to achieve within a single language. The proof assistant is used to design a specification, prove all the correctness invariants of interest and performs an extraction into a high-performance language of choice. This may take a non-trivial effort, but correctness *is* demanding!

Having a trusted specification as well as entire code-generation pipeline within a single dependently-typed framework is incredibly powerful. As we have demonstrated at the example of the neural network, we can introduce domain-specific optimisations and even swap the backend in case its performance is unsatisfying. For our example, the entire framework that includes array theory, DSL, optimisations and extraction is about 2000 lines of Agda code. We managed to introduce a new backend in about two days and match performance of the hand-written code.

A lot of pieces that we have developed in this paper can be reused in other numerical applications. We used dependent types to guarantee the absence of out-of-bound indexing, certain function being inverses as well as well-scopedness and well-typedness of our DSL. However, there are many more opportunities that we did not explore. For example, one can prove the correctness of optimisations,

<sup>5</sup>SaC compiler generates C code, so we can control what flags it uses when compiling it.

relating evaluation of optimised and non-optimised expressions. We can provide more guarantees when we run extraction, e.g. we can formalise some aspects of the backend language and relate them to our DSL. As for the DSL itself, we can try extending it with internal let constructions which should improve our C code generation as well as facilitate optimisations of the derivatives.

There are indeed plenty of opportunities, but the key point is this. Correctness and performance are competing requirements when it comes to application design. Therefore, such a cooperation between correctness-oriented and performance-oriented tools is likely to persist. With this work we demonstrate that such cooperation is feasible in practice.

## References

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. 2016. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*. 265–283.
- [2] Agda Development Team. 2024. Agda 2.6.3 documentation. <https://agda.readthedocs.io/en/v2.6.3/> Accessed [2024/02/28].
- [3] Atılım Günes Baydin, Barak A. Pearlmutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind. 2017. Automatic differentiation in machine learning: a survey. *J. Mach. Learn. Res.* 18, 1 (jan 2017), 5595–5637.
- [4] Aloïs Brunel, Damiano Mazza, and Michele Pagani. 2019. Backpropagation in the simply typed lambda-calculus with linear negation. *Proc. ACM Program. Lang.* 4, POPL, Article 64 (dec 2019), 27 pages. doi:10.1145/3371132
- [5] Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu, Chiyuan Zhang, and Zheng Zhang. 2015. MXNet: A Flexible and Efficient Machine Learning Library for Heterogeneous Distributed Systems. arXiv:1512.01274 [cs.DC]
- [6] Geoffrey S. H. Cruttwell, Bruno Gavranović, Neil Ghani, Paul Wilson, and Fabio Zanasi. 2022. Categorical Foundations of Gradient-Based Learning. In *Programming Languages and Systems*, Ilya Sergey (Ed.). Springer International Publishing, Cham, 1–28.
- [7] Conal Elliott. 2018. The simple essence of automatic differentiation. *Proc. ACM Program. Lang.* 2, ICFP, Article 70 (jul 2018), 29 pages. doi:10.1145/3236765
- [8] Brendan Fong, David Spivak, and Rémy Tuyéras. 2021. Backprop as functor: a compositional perspective on supervised learning. In *Proceedings of the 34th Annual ACM/IEEE Symposium on Logic in Computer Science (Vancouver, Canada) (LICS '19)*. IEEE Press, Article 11, 13 pages.
- [9] Clemens Grelck. 2005. Shared memory multiprocessor support for functional array processing in SAC. *Journal of Functional Programming* 15, 3 (2005), 353–401. doi:10.1017/S0956796805005538
- [10] James W. Hanson, Jane Shearin Caviness, and Camilla Joseph. 1962. Analytic differentiation by computer. *Commun. ACM* 5, 6 (jun 1962), 349–355. doi:10.1145/367766.368195
- [11] Chantal Keller and Thorsten Altenkirch. 2010. Hereditary substitutions for simple types, formalized. In *Proceedings of the Third ACM SIGPLAN Workshop on Mathematically Structured Functional Programming (Baltimore, Maryland, USA) (MSFP '10)*. Association for Computing Machinery, New York, NY, USA, 3–10. doi:10.1145/1863597.1863601
- [12] Faustyna Krawiec, Simon Peyton Jones, Neel Krishnaswami, Tom Ellis, Richard A. Eisenberg, and Andrew Fitzgibbon. 2022. Provably correct, asymptotically efficient, higher-order reverse-mode automatic differentiation. *Proc. ACM Program. Lang.* 6, POPL, Article 48 (jan 2022), 30 pages. doi:10.1145/3498710
- [13] Min Lin. 2024. Automatic Functional Differentiation in JAX. In *The Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=gzT61ziSCu>
- [14] Conor McBride and Ross Paterson. 2008. Applicative programming with effects. *Journal of Functional Programming* 18, 1 (2008), 1–13. doi:10.1017/S0956796807006326
- [15] J. F. Nolan. 1953. *Analytical Differentiation on a Digital Computer*. Master's thesis. Massachusetts Institute of Technology.
- [16] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. *PyTorch: an imperative style, high-performance deep learning library*. Curran Associates Inc., Red Hook, NY, USA.
- [17] F. Pfenning and C. Elliott. 1988. Higher-order abstract syntax. In *Proceedings of the ACM SIGPLAN 1988 Conference on Programming Language Design and Implementation (Atlanta, Georgia, USA) (PLDI '88)*. Association for Computing Machinery, New York, NY, USA, 199–208. doi:10.1145/53990.54010
- [18] R. Schenck, O. Rånnning, T. Henriksen, and C. E. Oancea. 2022. AD for an Array Language with Nested Parallelism. In *2022 SC22: International Conference for High Performance Computing, Networking, Storage and Analysis (SC) (SC)*. IEEE

Computer Society, Los Alamitos, CA, USA, 829–843. <https://doi.ieeecomputersociety.org/>

[19] Sven-Bodo Scholz. 1998. With-loop-folding in Sac — Condensing Consecutive Array Operations. In *Implementation of Functional Languages, 9th International Workshop (IFL'97), St. Andrews, UK, Selected Papers (Lecture Notes in Computer Science, Vol. 1467)*, Chris Clack, Tony Davie, and Kevin Hammond (Eds.). Springer, 72–92. [doi:10.1007/BFb0055425](https://doi.org/10.1007/BFb0055425)

[20] Sven-Bodo Scholz. 2003. Single Assignment C: efficient support for high-level array operations in a functional setting. *Journal of Functional Programming* 13, 6 (2003), 1005–1059. [doi:10.1017/S0956796802004458](https://doi.org/10.1017/S0956796802004458)

[21] Justin Slepak, Olin Shivers, and Panagiotis Manolios. 2014. An Array-Oriented Language with Static Rank Polymorphism. In *Proceedings of the 23rd European Symposium on Programming Languages and Systems - Volume 8410*. Springer-Verlag, Berlin, Heidelberg, 27–46. [doi:10.1007/978-3-642-54833-8\\_3](https://doi.org/10.1007/978-3-642-54833-8_3)

[22] Birthe van den Berg, Tom Schrijvers, James McKinna, and Alexander Vandenbroucke. 2024. Forward- or reverse-mode automatic differentiation: What’s the difference? *Science of Computer Programming* 231 (2024), 103010. [doi:10.1016/j.scico.2023.103010](https://doi.org/10.1016/j.scico.2023.103010)

[23] Artjoms Šinkarovs, Hans-Nikolai Vießmann, and Sven-Bodo Scholz. 2021. Array languages make neural networks fast. In *Proceedings of the 7th ACM SIGPLAN International Workshop on Libraries, Languages and Compilers for Array Programming (Virtual, Canada) (ARRAY 2021)*. Association for Computing Machinery, New York, NY, USA, 39–50. [doi:10.1145/3460944.3464312](https://doi.org/10.1145/3460944.3464312)