arXiv:1907.00509v1 [cs.PL] 1 Jul 2019

# *The Semantics of Rank Polymorphism*

Justin Slepak, Olin Shivers and Panagiotis Manolios
Northeastern University

## Abstract

Iverson's APL and its descendants (such as J, K and FISh) are examples of the family of "rank-polymorphic" programming languages. The principal control mechanism of such languages is the general lifting of functions that operate on arrays of rank (or dimension) $r$ to operate on arrays of any higher rank $r' > r$. We present a core, functional language, Remora, that captures this mechanism, and develop both a formal, dynamic semantics for the language, and an accompanying static, rank-polymorphic type system for the language. Critically, the static semantics captures the shape-based lifting mechanism of the language. We establish the usual progress and preservation properties for the type system, showing that it is sound, which means that "array shape" errors cannot occur at run time in a well-typed program. Our type system uses dependent types, including an existential type abstraction which permits programs to operate on arrays whose shape or rank is computed dynamically; however, it is restricted enough to permit static type checking.

The rank-polymorphic computational paradigm is unusual in that the types of arguments affect the dynamic execution of the program—they are what drive the rank-polymorphic distribution of a function across arrays of higher rank. To highlight this property, we additionally present a dynamic semantics for a *partially erased* variant of the fully-typed language and show that a computation performed with a fully-typed term stays in lock step with the computation performed with its partially erased term. The residual types thus precisely characterise the type information that is needed by the dynamic semantics, a property useful for the (eventual) construction of efficient compilers for rank-polymorphic languages.

## 1 Introduction

The essence of the rank-polymorphic programming model is implicitly treating all operations as *aggregate* operations, usable on arrays with arbitrarily many dimensions. The model was first introduced by Iverson with the language APL (Iverson, 1962). Over time, Iverson continued to develop this programming model, making it gradually more flexible, eventually leading to the creation of J (Jsoftware, Inc., n.d.) as a successor to APL. The boon APL offered programmers was a notation without loops or recursion: Programs would automatically follow a control-flow structure appropriate for the data being consumed. The nature of the implicit iteration structure could be modified using second-order operators, such as folding, scanning, or operating over a moving window. These second-order operators would directly reveal all loop-carried data dependences.

In this sense, other languages demanded that unnecessary work be put into both compilers and user programs. The programmer would be expected to write the program's iteration structure explicitly; in many languages this entails describing a particular serial encoding

of what is fundamentally parallelizable computation. The compiler must then perform intricate static analysis to see past the programmer's overspecified iteration schedule.

The design of APL earned a Turing award for Iverson (Iverson, 1980) as well as a mention in an earlier Turing lecture (Backus, 1978), praising it for showing the basis of a solution to the "ven Neumann bottleneck." However APL's subsequent development proceded largely in isolation from mainstream programming-language research. The APL family of languages painted itself into a corner with design decisions such as requiring functions to take only one or two arguments and making parsing dependent on values assigned at run time. As a result, APL compilers were forced to support only a subset of the language (such as Budd's compiler (Budd, 1988)) or to operate on small sections of code, alternating between executing each line of the program and compiling the next one (Johnston, 1979). What we gain from the rank-polymorphic programming model's natural friendliness to parallelism, we can easily lose by continually interrupting the program to return control to a line-at-a-time compiler. Limiting the compiler to operating over a narrow window of code can also eliminate opportunities for code transformations like fusion, forcing unnecessary materialization of large arrays.

The tragedy of rank-polymorphic programming does not end at forgone opportunities for performance. Despite the convenience of rank polymorphism for writing array-processing code—a common task in many application domains—APL and its close descendants do not see widespread use. There is enough desire for implicitly aggregate computation to support user communities for systems such as NumPy (Oliphant, 2006) and MATLAB (Mathworks, 1992), which do not follow as principled or as flexible a rule for matching functions with aggregate arguments[1]. However, programmers are driven away from APL itself by features such as obtuse syntax, restrictions on function arity, poor support for naming things, and a limited universe of atomic data to populate the arrays (Abrams, 1975).

Our goal is to study rank polymorphism itself without getting bogged down by APL's other baggage. A formal semantics of rank polymorphism is the essential groundwork for understanding how rank-polymorphic programs ought to behave, how they should be compiled, and how they can be safely transformed to reduce execution cost. To that end, we develop Remora, a language which integrates rank polymorphism with typed $\lambda$-calculus.

A key problem impeding static compilation of rank-polymorphic programs is identifying the implicit iteration structure at each function application. Even without obscuring the programming model with the idiosyncratic special case behavior APL accreted, rank polymorphism itself has seemed "too dynamic" for good static compilation due to having its control structure derived from computed data. The old style of line-at-a-time compilation relied on inspecting functions and data at run time to decide what loop structure to emit.

Remora's answer to the problem of finding the iteration space is a type system which describes the shapes of arrays and thereby identifies the implicit iteration space for each function application. In order for types to provide enough detail about array shapes, we use a restricted form of dependent typing, in the style of Dependent ML (Xi, 1998). In Dependent ML, types are not parameterized over arbitrary program terms but over a much

---

[1] For example, operations which already expect aggregate data—for example if the programmer writes a function to compute the norm of a vector or the determinant of a matrix—do not always lift easily to consume even higher-dimensional arguments

more restricted language. For Remora, our language of type indices consists of natural numbers, describing individual dimensions, and sequences of natural numbers, describing array shapes.

Past work on applying dependent types to computing with arrays has focused on ensuring the safety of accessing individual array elements (Xi & Pfenning, 1998; Trojahner & Grelck, 2009). Bounds checking array indices is essential in a programming model where extracting a single element is the only elimination form for arrays, but the rank-polymorphic programming model generally eschews this operation. Instead, arrays are consumed whole, and function application itself serves as the elimination form for arrays.

In Remora, a function's type describes the shapes of the expected argument arrays, called the "cells," and the type of the atomic data inside the array. The typing rule for function application is responsible for identifying the "frame" shape, *i.e.*, the iteration structure derived from the non-cell dimensions. Type soundness means that the type system produces more than a safety guarantee: conclusions it draws about the iteration structure can be used to correctly compile the program. Our type system is flexible enough to express polymorphism over the cell shape, such as a determinant function that can operate on square matrix cells of any size. It can also handle functions whose output shape is not determined by input shape alone, such as reading a vector of unknown size from user input or generating an array of caller-specified shape.

We begin with an overview of the rank-polymorphic programming model, written as a programming tutorial for an untyped variant of Remora. After developing the intuition for rank polymorphism, we present a formal description of Remora's core language. This includes Remora's abstract syntax, the language of type indices with its associated theory, the static semantics which identifies array shapes and iteration spaces, a type-driven dynamic semantics, and a type-soundness theorem linking the static and dynamic semantics. Since our formal presentation is intrinsically typed, we also include an algorithm for partial type erasure, to characterize which type-level information is truly necessary to keep at run time. A bisimulation argument connects the dynamic semantics of explicitly-typed Remora to that of erased Remora.

## 2 Formalism

We present a formal description of Core Remora, which describes the control-flow mechanism used for computing on arrays. While the nested-vector shorthand used earlier is convenient for human use, this formalism explicitly distinguishes atoms from arrays.

The basic design goal for the language of types and indices is to describe the program's control structure, for a compiler's benefit. This requires a detailed description of array shapes. Knowing only the number of axes an array has is insufficient because good mapping from source to hardware—*e.g.*, whether to emit vector instructions, invoke a GPGPU kernel, or fork separate parallel threads—depends more on the actual sizes of individual axes than on how many there are. We use indexed types, in the style of Dependent ML: rather than allowing types to be parameterized over arbitrary terms, they are parameterized over a limited language of type indices. Remora's index language consists of natural numbers, representing individual dimensions, and sequences of naturals, representing array shapes (or fragments of shapes). So the type of an array has the form (Arr $\tau$ $\iota$), where

$\tau$ identifies the type of the array's atoms, and $\iota$ describes the array's shape. This includes enough detail for the type system to describe how the function and argument arrays align in function application. It also grants the ability to statically detect arrays that cannot be properly aligned.

Fixed-size computation, requiring every function to exactly specify its argument and result sizes, is far too restrictive for practical use. Programmers should not have to write a separate `vector-mean` function for every possible length of vector their programs might use. So the index language must permit variables, and the type language must allow universal quantification. This is phrased as a dependent product: `(Pi ((x γ) ... ) τ)`. Each $x$ is marked with its sort $\gamma$, which specifies whether $x$ ranges over individual dimensions ($\gamma = $ `Dim`) or sequences ($\gamma = $ `Shape`). Now we can give a type to `vector-mean`:

```
(Pi ((n Dim))
  (-> ((Arr Float (Shp n)))
      (Arr Float (Shp))))
```

This function will lift to operate on higher-rank arrays of `Float`s, effectively behaving as a minor-axis mean function. Having `Shape` variables in addition to `Dim` variables allows us to type a major-axis mean function as well:

```
(Pi ((c Shape) (n Dim))
  (-> ((Arr Float (++ (Shp n) c)))
      (Arr Float c)))
```

Combined with parametric polymorphism, where type variables can be quantified separately over the kinds `Atom` and `Array`, we now have a lot of flexibility in describing a function's behavior. For example, `append` stitches two arrays together along their major axis. This requires that the $(n-1)$-dimensional pieces of each $n$-dimensional array have the same type (*i.e.*, they must share the same `Atom`-kinded type variable), but we must introduce separate index variables (of sort `Dim`) for the arguments' major axes. The type we give to `append` is

```
(Pi ((c Shape) (m Dim) (n Dim))
  (Forall ((a Atom))
    (-> ((Arr a (++ (Shp m) c))
         (Arr a (++ (Shp n) c)))
        (Arr a (++ (Shp (+ m n)) c)))))
```

We have index variables `m` and `n` to stand for each argument's first dimension and `c` to denote the rest of their shapes. Quantifying over the type variable `a` allows `append` to work independent of the type of atoms its arguments contain. Our result type's first dimension is the sum of the arguments' first dimensions, but it has the same "remainder" shape `c`.

Quantifying over `Array` types is a convenience—it is not strictly necessary. Concrete `Array`-kinded types must be of the form `(Arr τ ι)`, so polymorphic types of the form `(Forall ((x Array)) τ)` could be rewritten with fresh variables $x_s$ and $x_a$ as

```
(Pi ((xₛ Shape))
  (Forall ((xₐ Atom))
    τ[x ↦ (Arr xₐ xₛ)]))
```

Even having escaped the confines of fixed-size computation, we so far only have functions whose result shape depends solely on its arguments' shapes. Common utility functions such as `iota` and `filter` have result shapes which depend on the actual run-time data they receive. We can solve this limitation using existential quantification. A dependent sum type, (Sigma ($(x\ \gamma)\ \ldots$) $\tau$), conceptually represents a tuple containing indices whose respective sorts are $\gamma \ldots$ and an array whose type $\tau$ may depend on those indices. This is the type-level description of a box, the atomic wrapper around an arbitrary array. Such types can encode arrays whose dimensions are not all known. For example, (Sigma ((n Dim)) (Arr Int (Shp n))) can describe any vector of integers without giving its specific length. Note that the type still requires the underlying array to have rank 1. No scalar or matrix or higher-ranked array can fit the pattern specified by (Arr Int (Shp n)), the dependent sum's body. Use of dependent sums offers a lot of freedom in stating what shape information is known precisely and what is hidden. At one extreme, (Sigma ((s Shape)) (Arr Int s)) could contain an array of absolutely any shape. We can also write more detailed descriptions, *e.g.*, floating-point matrices with exactly three rows and at least two columns:

```
(Sigma ((c Dim))
  (Arr Float
       (Shp 3 (+ 2 c)))))
```

Quantifying over shapes allows us to describe arrays where only specific axes are known, such as boolean arrays whose leading axis has length 10:

```
(Sigma ((s Shape))
  (Arr Bool
       (++ (Shp 10) s)))
```

Typing boxes as dependent sums also permits controlled access to ragged arrays, which are typed as arrays of boxes. Consider a vector of 20 strings of varying lengths:

```
(Arr (Sigma ((len Dim))
       (Arr Char (Shp len)))
     (Shp 20))
```

Any function written to operate on a box of the appropriate type, in this case containing `Char` vectors of completely unknown length, can be safely lifted to operate on this vector of strings. Separating the lifting over the outer dimensions from the lifting over inner, existentially hidden dimensions reflects an important consideration for code generation: ragged dimensions in a type identify when implicit parallelism is *irregular*, in contrast with the strictly regular parallelism offered in box-free code. Raggedness is not restricted to the minor axis because a box's type can still specify some exact dimensions, as in

```
(Arr (Sigma ((l Dim))
       (Arr Char (Shp l 80)))
     (Shp 20))
```

Here we have a vector of 20 documents, each of which is a character array containing an unknown number of 80-character lines.

| | | |
|---|---|---|
| $e \in Expr$ ::= | | *Expressions* |
| | $x$ | *Variable reference* |
| &#124; | $(\texttt{array}\,(n \dots)\,\mathfrak{a} \dots)$ | *Array, containing atoms* |
| &#124; | $(\texttt{array}\,(n \dots)\,\tau)$ | *Empty array, with its atom type* |
| &#124; | $(\texttt{frame}\,(n \dots)\,e \dots)$ | *Frame, containing array cells* |
| &#124; | $(\texttt{frame}\,(n \dots)\,e \dots)$ | *Empty frame, with its cell type* |
| &#124; | $(e_f\,e_a \dots)$ | *Term application* |
| &#124; | $(\texttt{t-app}\,e\,\tau \dots)$ | *Type application* |
| &#124; | $(\texttt{i-app}\,e\,\iota \dots)$ | *Index application* |
| &#124; | $(\texttt{unbox}\,(x_i \dots x_e\,e_s)\,e_b)$ | *Let-binding box contents* |
| $v \in Val$ ::= $x$ &#124; $(\texttt{array}\,(n \dots)\,\mathfrak{v} \dots)$ | | *Values* |
| $\mathfrak{a} \in Atom$ ::= | | *Atoms* |
| | $\mathfrak{b}$ | *Base value* |
| &#124; | $\mathfrak{o}$ | *Primitive operator* |
| &#124; | $(\lambda\,((x\,\tau) \dots)\,e)$ | *Term abstraction* |
| &#124; | $(\texttt{T}\lambda\,((x\,k) \dots)\,v)$ | *Type abstraction* |
| &#124; | $(\texttt{I}\lambda\,((x\,\gamma) \dots)\,v)$ | *Index abstraction* |
| &#124; | $(\texttt{box}\,\iota \dots e\,\tau)$ | *Boxed array* |
| $\mathfrak{v} \in Atval$ ::= $\mathfrak{b}$ &#124; $\mathfrak{o}$ &#124; $(\lambda\,((x\,\tau) \dots)\,e)$ &#124; $(\texttt{T}\lambda\,((x\,k) \dots)\,v)$ | | *Atomic values* |
| &#124; | $(\texttt{I}\lambda\,((x\,\gamma) \dots)\,v)$ &#124; $(\texttt{box}\,\iota \dots v\,\tau)$ | |
| $\tau \in Type$ ::= | | *Types* |
| | $x$ | *Type variable* |
| &#124; | $B$ | *Base type* |
| &#124; | $(\texttt{Arr}\,\tau\,\iota)$ | *Array* |
| &#124; | $(\texttt{->}\,(\tau \dots)\,\tau')$ | *Function* |
| &#124; | $(\texttt{Forall}\,((x\,k) \dots)\,\tau)$ | *Universal* |
| &#124; | $(\texttt{Pi}\,((x\,\gamma) \dots)\,\tau)$ | *Dependent product* |
| &#124; | $(\texttt{Sigma}\,((x\,\gamma) \dots)\,\tau)$ | *Dependent sum* |
| $k \in Kind$ ::= $\texttt{Array}$ &#124; $\texttt{Atom}$ | | *Kinds* |
| $\iota \in Idx$ ::= | | *Type indices* |
| | $x$ | *Type variable* |
| &#124; | $n$ | *Single dimension* |
| &#124; | $(\texttt{Shp}\,\iota \dots)$ | *Sequence of dimensions)* |
| &#124; | $(\texttt{+}\,\iota \dots)$ | *Adding dimensions* |
| &#124; | $(\texttt{++}\,\iota \dots)$ | *Appending shapes* |
| $\gamma \in Sort$ ::= $\texttt{Shape}$ &#124; $\texttt{Dim}$ | | *Index sorts* |
| $\mathfrak{o} \in Op$ ::= $\texttt{+}$ &#124; $\texttt{-}$ &#124; $\texttt{*}$ &#124; $\texttt{/}$ &#124; $\texttt{append}$ &#124; $\texttt{reduce}$ &#124; $\texttt{iota}$ &#124; $\dots$ | | *Primitive operators* |
| $\mathfrak{f} \in Func$ ::= $\mathfrak{o}$ &#124; $e$ | | *Functions* |
| $t \in Term$ ::= $\mathfrak{a}$ &#124; $(\lambda\,((x\,\tau) \dots)\,e)$ | | *Terms* |

Fig. 1. Core Remora grammar

### *2.1 Syntax*

The grammar for Core Remora is given in Figure 1. Term-level syntax is divided into
atoms, noted as $\mathfrak{a}$, and expressions, noted as $e$. Expressions produce arrays, which contain
atoms. For the most part, atom terms perform only trivial computation. This rule applies
to base values, noted as $\mathfrak{b}$; primitive operators, noted as $\mathfrak{o}$; and $\lambda$-abstractions, which
may abstract over terms, types, and type indices. As an exception, a box gives an atomic
view of an array of any shape and may therefore perform any computation to compute its
contents. A box hides part of its contents' shape, using a dependent sum. It existentially

quantifies type indices, but an explicit type annotation is required. A box built from the index 3 and a $3 \times 3$ matrix could be meant, for example, as an unspecified-length vector containing 3-vectors, with type `(Sigma ((n Dim)) (Arr Int (Shp n 3)))` or as a square matrix of unspecified size, with type `(Sigma ((n Dim)) (Arr Int (Shp n n))) (Sigma ((n Dim)) (Arr Int (Shp n n)))`;.

An array can be written as a literal, with its shape and individual atoms listed directly. It can also be written in nested form as a frame containing cells (its subexpressions) arranged in the specified shape. For example, the matrix $\left[\begin{smallmatrix} 1 & 2 \\ 3 & 4 \end{smallmatrix}\right]$ can be written as the literal

```
(array (2 2) 1 2 3 4)
```

or as a vector frame of vector literal cells:

```
(frame (2) (array (2) 1 2) (array (2) 3 4))
```

The frame notation allows construction of arrays from unevaluated cells. An empty array (*i.e.*, one with a zero in its shape) must be written with the type its elements are meant to have. An empty vector of integers is a different value than an empty vector of booleans, and they inhabit different types.

Term, type, and index abstractions can be applied to zero or more expressions, types, or indices. The body of the abstraction must itself be an expression, *i.e.*, all functions produce arrays as their results.

Consuming a box let-binds its index- and term-level contents. Suppose we have M, a boxed square matrix of unspecified size. Unboxing M as in `(unbox (l a M) e)` lets us use the index variable `l` and term variable `a` within *e*, the body.

Types include base types (noted as *B*), functions, arrays, universal types, and dependent products and sums. Universals specify the kind of each type argument, and dependent products and sums specify the sort of each index argument. Types are classified as either `Atom` or `Array`. Type indices are naturals and sequences of naturals, with addition and appending as the only operators. They are classified into sorts, `Dim` and `Shape`.

The grammar in Figure 1 does not require any specific set of primitive operators, base types, and base values. An example collection of array-manipulation primitives and their types is given in Figure 2. For readability, we elide the enclosing `Pi` and `Forall` forms. Most of these primitives perform some operation along the argument's major axis. For example, `head` extracts the first scalar of a vector, the first row of a matrix, *etc.* This means that the argument shape must have one dimension more than the result shape, and that extra dimension must be nonzero. This is expressed in the type of `head` by giving the argument shape `(++ (+ 1 d) s)`, *i.e.*, a single dimension which is 1 plus any arbitrary natural followed by any arbitrary sequence of naturals. In taking one scalar from a 3-vector, we would instantiate `d` as the dimension 2 and `s` as the empty shape `(Shp)`. If we want to extract the first plane of a $5 \times 6 \times 7$ array, we use 4 for `d` and `(Shp 6 7)` for `s`.

Since these operations work along the major axis, we can use other axes instead by instantiating them differently. Suppose `mtx` is the matrix `(array (3 2) 0 1 2 3 4 5)`, which has type `(Arr Num (Shp 3 2))`. Then `(t-app (i-app head 2 (Shp 2)) Num)` is a function which extracts the first row of a $(1+2) \times 2$ (*i.e.*, $3 \times 2$) matrix. So `((t-app (i-app head 2 (Shp 2)) Num) mtx)` evaluates to `(array (2) 0 1)`. Instead, consider `(t-app (i-app head 1 (Shp)) Num)`. This is a function with input type `(Arr`

`Num (Shp 2))` and output type `(Arr Num (Shp))`. It extracts the first scalar of a 2-vector. When applied to `mtx`, this function *lifts* to extract the first scalar from *each* 2-vector, gathering the results as `(frame (3) (array () 0) (array () 2) (array () 4))`. Evaluation proceeds, reducing this to `(array (3) 0 2 4)`, the first column of `mtx`.

Several primitives must return boxed arrays because the type system cannot keep track of enough information to fully describe the result shape. As an extreme example, `read-nums` reads a vector of numbers from user input, and there is no way of knowing until run time how long a vector the user will enter. In other cases, the necessity of boxing comes from a limit on the type system's expressive power. The `ravel` function produces a vector whose atoms are all those of the argument array, laid out in row-major order. The length of the `ravel` of some array is fully determined by that array's shape: it is the product of all of its dimensions. However the undecidability of Peano arithmetic would interfere with type checking (not to mention future efforts on type inference). Since "product of all dimensions" is not expressible in Presburger arithmetic, we instead have `ravel` return a boxed vector.

Boxing is not limited to vectors. For example, `filter` uses a vector of booleans to decide which parts of an array to retain. Since the number of true entries in that vector is unknown, the size of the result's major axis is also unknown. The resulting `Sigma` type existentially quantifies only that one dimension, and leaves the remaining dimensions externally visible.

The `iota` functions and their variants , described in Figure 3, form a useful case study on what invariants can be expressed in Remora's type system. These functions produce arrays whose atoms are successive natural numbers starting from 0, such as `(array (2 3) 0 1 2 3 4 5)`, representing the matrix $\begin{bmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \end{bmatrix}$. The argument to `iota` is a vector of numbers specifying the result array's shape. Since this vector can be dynamically computed, we cannot give any specific shape for `iota`'s return type. Instead, `iota` must return a box with existentially quantified shape. Recall that boxing arrays allows functions with data-dependent result shape to lift safely, since applying `iota` to `(array (2 2) 3 3 4 4)` must produce a $3 \times 3$ matrix and a $4 \times 4$ matrix as its two result cells.

Variants on `iota` allow the programmer to communicate more detailed knowledge to the type system. When the result is meant to be a vector, `iota/v` takes that vector's length as the argument. The resulting box is typed as a vector of unknown length rather than an array of completely unknown shape. Knowing that we have a vector of numbers rather than any arbitrary array means, for example, that summing the box's contents with `reduce` is certain to produce a scalar. We can therefore type the following function as consuming and producing *non-boxed* scalar numbers:

```
(λ ((n (Arr Num (Shp))))
  (unbox (len nums ((array () iota/v) n))
    ((t-app (i-app (array () reduce) len (Shp)) Num)
     +
     ((t-app (i-app (array () append) 1 len (Shp)) Num)
      (array () 0)
      nums))))
```

| Function | Type |
|---|---|
| head, tail | `(-> ((Arr t (++ (Shp (+ 1 d)) s)))`<br>`    (Arr t s))` |
| behead, curtail | `(-> ((Arr t (++ (Shp (+ 1 d)) s)))`<br>`    (Arr t (++ (Shp d) s)))` |
| length | `(-> ((Arr t (++ (Shp d) s)))`<br>`    (Arr Num (Shp)))` |
| shape, ravel | `(-> ((Arr t s))`<br>`    (Arr (Sigma ((d Dim)) (Arr Num (Shp d)))`<br>`         (Shp)))` |
| append | `(-> ((Arr t (++ (Shp m) s))`<br>`     (Arr t (++ (Shp n) s)))`<br>`    (Arr t (++ (Shp (+ m n)) s)))` |
| reverse | `(-> ((Arr t (++ (Shp d) s)))`<br>`    (Arr t (++ (Shp d) s)))` |
| rotate | `(-> ((Arr t (++ (Shp d) s))`<br>`     (Arr Num (Shp)))`<br>`    (Arr t (++ (Shp d) s)))` |
| fold | `(-> ((Arr (-> ((Arr t s) T) T) (Shp))`<br>`     T`<br>`     (Arr t (++ (Shp d) s)))`<br>`    T)` |
| reduce | `(-> ((Arr (-> ((Arr t s) (Arr t s))`<br>`               (Arr t s))`<br>`         (Shp))`<br>`     (Arr t (++ (Shp (+ 1 d)) s)))`<br>`    (Arr t s))` |
| scan | `(-> ((Arr (-> ((Arr u r) (Arr t s)) (Arr u r)) (Shp))`<br>`     (Arr u r)`<br>`     (Arr t (++ (Shp d) s)))`<br>`    (Arr u (++ (Shp d) r)))` |
| filter | `(-> ((Arr Bool d)`<br>`     (Arr t (++ (Shp d) s)))`<br>`    (Arr (Sigma ((k Dim)) (Arr t (++ (Shp k) s))) (Shp)))` |
| read-nums | `(-> () (Arr (Sigma ((k Dim)) (Arr Num (Shp k))) (Shp)))` |
| iota | `(-> ((Arr Num (Shp d)))`<br>`    (Arr (Sigma ((s Shape)) (Arr Num s)) (Shp)))` |
| reshape | `(-> ((Arr Num (Shp d))`<br>`     (Arr t r))`<br>`    (Arr (Sigma ((s Shape)) (Arr Num s)) (Shp)))` |

Fig. 2. Common array-manipulation primitive operations and their Remora types

| Function | Type |
|---|---|
| iota | `(-> ((Arr Num (Shp d)))` `(Arr (Sigma ((s Shape)) (Arr Num s)) (Shp)))` |
| iota/v | `(-> ((Arr Num (Shp)))` `(Arr (Sigma ((d Dim)) (Arr t (Shp d))) (Shp)))` |
| iota/s | `(Pi ((s Shape))` `(-> () (Arr Num s)))` |
| iota/w | `(-> ((Arr t s))` `(Arr Num s))` |

Fig. 3. Types for `iota` and its variants

In a more programmer-friendly surface language, with automatic instantiation of polymorphic functions[2] and conversion of bare atoms to scalar arrays, this might be written as:

```
(λ ((n (Arr Num (Shp))))
  (unbox (len nums (iota/v n))
    (reduce + (append [0] nums))))
```

Alternatively, the programmer might prefer to use `iota/s` to pass the desired result shape as a type index rather than as a term-level vector. In that case, there is no need to box the result array. In the automatic-instantiation shorthand, `iota/s` may be stylistically awkward, calling for the variant `iota/w`, which takes an extra array argument as a "shape witness" rather than instantiating at a shape index. Producing a number array whose shape matches some existing array `xs` could then be written as `(iota/w xs)` instead of `((i-app iota/s shape-of-xs))`.

The `reshape` function behaves similarly to `iota`, except that the atoms in the result array are drawn from the second argument, repeating them cyclically if necessary. So using `reshape` with the shape specification `(array (2) 3 2)` and the vector `(array (5) 1 2 3 4 5)` produces the $3 \times 2$ matrix `(array (3 2) 1 2 3 4 5 1)`. Like `iota`, `reshape` benefits from alternative ways for the programmer to specify the result shape.

### 2.2  Theory of type indices

Type indices, given in program syntax as $\iota$, represent individual dimensions, taken from $\mathbb{N}$, and array shapes, taken from the free monoid on $\mathbb{N}$. The theory of the free monoid on $\mathbb{N}$ includes as axioms the associativity of adding naturals and appending sequences as well as unique identity elements for addition (zero) and appending (the scalar shape, $\square$):

$$0 + i = i + 0 = i$$

$$(i + j) + k = i + (j + k)$$

$$\square \mathbin{+\!\!\!+} a = a \mathbin{+\!\!\!+} \square = a$$

---

[2]  This inference problem is beyond the scope of this paper.

Fig. 4. Overlap axiom, visualized: $w$ is the overlapping portion of $a$ and $d$.

$$(a +\!\!\!+ b) +\!\!\!+ c = a +\!\!\!+ (b +\!\!\!+ c)$$

As the *free* monoid, it also follows an equidivisibility rule which states that if two uses of the append operator give the same result, there is some completing subsequence, representing the overlap between each use's larger argument (demonstrated in Figure 4):

$$a +\!\!\!+ b = c +\!\!\!+ d \implies \exists w.\,(a +\!\!\!+ w = c \wedge w +\!\!\!+ d = b) \vee (c +\!\!\!+ w = a \wedge w +\!\!\!+ b = d)$$

A free monoid (on any set of generators) also has a homomorphism to the monoid formed by $\mathbb{N}$ under addition, with the property that only the free monoid's identity element can be mapped to 0. This can be axiomatized with one additional function symbol $L$:

$$L(a) = 0 \implies a = \square$$

$$L(a +\!\!\!+ b) = L(a) + L(b)$$

Using equidivisibility and the homomorphism to the additive $\mathbb{N}$ monoid, we can define a partial operator $\dot{-}$ for prefix subtraction: $a \dot{-} b = c$ iff $b +\!\!\!+ c = a$. For example, $[3,4,5,6] \dot{-} [3,4] = [5,6]$, whereas $[3,4,5,6] \dot{-} [4]$ is undefined.

Type checking only requires a very restricted fragment of this theory. Pairs of indices are only checked for equality in isolation from each other, and no information about an index (other than its sort) is given in the program. So the check is for the validity of a single equality—no connectives or quantifiers needed. This fragment can be decided efficiently by comparing indices written in canonical form. Two `Dims` which are equal must simplify to sums with the same constant component and the same coefficient on corresponding variables. For example,

$$(+\, x\ y\ 5\ x) = (+\ (+\, x\ x)\ 5\ y)$$

is valid because both simplify to $2x + y + 5$, whereas

$$(+\, q\ 5\ y) = (+\ (+\, x\ x)\ 5\ y)$$

is false for any interpretation which does not assign $q$ to twice the value assigned to $x$ (and thus is not valid).

To decide the validity of an equality on `Shapes` (*i.e.*, sequences of naturals), we can again test by conversion to a canonical form: a sequence is written out as the concatenation of single `Dims` and `Shape` variables. Sorting rules guarantee that the individual elements of a sequence are natural numbers, and associativity permits nested appends to be collapsed away. Thus the index

$$(+\!\!+\ (\text{Shp}\ 2\ (+\, x\ 5\ x))\ (+\!\!+\ \ d\ (\text{Shp}\ 3)))$$

canonicalizes to

$$\texttt{(++ (Shp 2) (Shp (+ x x 5)) d (Shp 3))}$$

To show that this process does produce a canonical form, consider two shapes in this form which differ, and focus on the leftmost differing position in their respective lists of appended components. If they are syntactically different singleton shapes—their respective contents are two different canonicalized naturals—then an assignment under which those naturals differ will also make the full shapes differ at this position. If one is a singleton (`Shp` $\iota$) and the other a variable `s` (of sort `Shape`), then an interpretation which assigns the variables in $\iota$ such that its components sum to $n$ may also assign `s` to be the shape (`Shp (+ n 1)`). Again, an interpretation forces the shapes to be unequal. Finally, if this position has variables `s` and `t`, choose an interpretation mapping `s` to (`Shp 1`) and `t` to (`Shp 2`) to produce unequal interpretations of the whole shapes.

Although type checking itself only requires this canonicalization process, constraint-based type inference would call for a more sophisticated solver due to the use of existential variables for choosing pieces of an index.

In order to describe when arguments' shapes are compatible, it is useful to impose a lattice structure on the universe of shapes. The lattice is built with the order $\sqsubseteq$ meaning that one shape is a prefix of another; a $\top$ is added to represent the join of incompatible shapes (we already have $\bot = \square$, as the empty shape is a prefix of every shape). For shapes $s_0$ and $s_1$, we have $s_0 \sqcup s_1 \neq \top$ if and only if $s_0 \sqsubseteq s_1$ or $s_1 \sqsubseteq s_0$. Generalizing to arbitrary finite joins, $\bigsqcup\{s \ldots\} \neq \top$ implies that the shapes $s \ldots$ are totally ordered, and the lattice structure means the shapes' join is one of the shapes themselves.

### 2.3 Static Semantics

Typing Core Remora uses a three-part environment structure: $\Theta$ is a partial function mapping index variables to sorts; $\Delta$ maps type variables to their kinds; and $\Gamma$ maps term variables to their types. The stratification of Dependent ML-style types allows indices to be checked using only the sort environment and types using only the sort and kind environments. Following the definition of each judgment form, we give a handful of lemmas which will be needed for a type soundness argument in Subsection 2.5. The well-formedness judgments each come with a lemma stating that the judgment gives a unique result to each well-formed term and that unique result is preserved by substituting well-formed assignments for free variables. When we show type soundness for Remora, these results will be needed to prove the preservation lemma. Uniqueness of typing is particularly important for Remora, where the implicit iteration in function application (including index and type abstractions) is driven by the types ascribed to the function and argument expressions. Well-defined program behavior relies on having a unique decomposition of each array into a frame of cells.

#### 2.3.1 Sorting

Figure 5 defines the sorting judgment, $\Theta \vdash \iota :: \gamma$, which states that in sort environment $\Theta$, the index $\iota$ has sort $\gamma$. Natural number literals have sort `Dim`. A sequence of indices is a `Shape`, provided that every element of the sequence is a `Dim`. Addition is used on `Dim`

$$\boxed{\Theta \vdash \iota :: \gamma}$$

$$\frac{n \in \mathbb{N}}{\Theta \vdash n :: \mathtt{Dim}} \text{ S-Nat} \qquad \frac{(x :: \gamma) \in \Theta}{\Theta \vdash x :: \gamma} \text{ S-Var} \qquad \frac{\Theta \vdash \iota_j :: \mathtt{Dim} \text{ for each } j}{\Theta \vdash (\mathtt{Shp}\ \iota \dots) :: \mathtt{Shape}} \text{ S-Shape}$$

$$\frac{\Theta \vdash \iota_j :: \mathtt{Dim} \text{ for each } j}{\Theta \vdash (+\ \iota \dots) :: \mathtt{Dim}} \text{ S-Plus} \qquad \frac{\Theta \vdash \iota_j :: \mathtt{Shape} \text{ for each } j}{\Theta \vdash (\mathtt{++}\ \iota \dots) :: \mathtt{Shape}} \text{ S-Append}$$

Fig. 5. Sorting rules

arguments to produce a `Dim`. `Shape` arguments may be appended, to form another `Shape`. Variables may be bound at either sort, but they can only be introduced into the environment by index abstraction and unboxing terms—the index language itself has no binding forms.

We give two results about the well-behaved nature of the sorting rules: No index inhabits both sorts (in the same environment), and replacing an index's variables with appropriately-sorted indices does not change the sort.

*Lemma 2.1* (*Uniqueness of sorting*)
If $\Theta \vdash \iota :: \gamma$ and $\Theta \vdash \iota :: \gamma'$, then $\gamma = \gamma'$.

*Proof*
No non-variable index form is compatible with multiple sorting rules, so they can only have whichever sort their one compatible rule concludes. It remains to show that uniqueness holds for variables. Since $\Theta$ is a well-defined partial function, mapping variables to sorts, $\Theta(x)$ can only have one value. If $\Theta(x) = \gamma$ and $\Theta(x) = \gamma'$, $\gamma = \gamma'$. $\qquad \square$

*Lemma 2.2* (*Preservation of sorts under index substitution*)
If $\Theta, x :: \gamma_x \vdash \iota :: \gamma$ and $\Theta \vdash \iota_x :: \gamma_x$ then $\Theta \vdash \iota[x \mapsto \iota_x] :: \gamma$.

*Proof*
This is straightforward induction on the original sort derivation. $\qquad \square$

### 2.3.2 Kinding

Kinding rules are given in figure 6. The `Array` kind is only ascribed to types built by the array type constructor and type variables bound at that kind. The array type constructor requires as its arguments an `Atom` type and a `Shape` index. Base types are fundamental, non-aggregate types, such as `Float` or `Bool`, so they are `Atoms`. Function types have kind `Atom`, but their input and output types must be `Arrays`. This reflects the rule that application is performed on arrays, and the function produces an array result. Similarly, universal types and dependent products, describing type and index abstractions, must have an `Array` as their body, while they themselves are `Atoms`. This rules out types whose inhabitants would have to be syntactically illegal due to containing expressions instead of atoms as their bodies. Since boxes present arrays as atoms, dependent sum types also have an `Array` body and are kinded as `Atoms`. A universal type adds bindings for its quantified type variables to $\Delta$. Dependent products and sums do the same for their index variables in $\Theta$.

As with sorting of indices, we expect a well-kinded type to inhabit only a single kind (fixing a particular environment). The kinding system should also allow free index or

$$\boxed{\Theta;\Delta \vdash \tau :: k}$$

$$\frac{(x :: k) \in \Delta}{\Theta;\Delta \vdash x :: k} \text{ K-Var} \qquad \frac{}{\Theta;\Delta \vdash B :: \texttt{Atom}} \text{ K-Base} \qquad \frac{\begin{array}{c}\Theta;\Delta \vdash \tau_j :: \texttt{Array} \text{ for each } j \\ \Theta;\Delta \vdash \tau' :: \texttt{Array}\end{array}}{\Theta;\Delta \vdash (\texttt{->} (\tau \dots) \ \tau') :: \texttt{Atom}} \text{ K-Fun}$$

$$\frac{\Theta;\Delta, x :: k \dots \vdash \tau :: \texttt{Array}}{\Theta;\Delta \vdash (\texttt{Forall} ((x \ k) \dots) \ \tau) :: \texttt{Atom}} \text{ K-Univ} \qquad \frac{\Theta, x :: \gamma \dots ;\Delta \vdash \tau :: \texttt{Array}}{\Theta;\Delta \vdash (\texttt{Pi} ((x \ \gamma) \dots) \ \tau) :: \texttt{Atom}} \text{ K-Pi}$$

$$\frac{\Theta, x :: \gamma \dots ;\Delta \vdash \tau :: \texttt{Array}}{\Theta;\Delta \vdash (\texttt{Sigma} ((x \ \gamma) \dots) \ \tau) :: \texttt{Atom}} \text{ K-Sigma}$$

$$\frac{\Theta \vdash \iota :: \texttt{Shape} \qquad \Theta;\Delta \vdash \tau :: \texttt{Atom}}{\Theta;\Delta \vdash (\texttt{Arr} \ \tau \ \iota) :: \texttt{Array}} \text{ K-Array} \qquad \text{K-Array}$$

Fig. 6. Kinding rules

type variables to be replaced with appropriately sorted or kinded indices or types without changing the original type's kind.

*Lemma 2.3 (Uniqueness of kinding)*
If $\Theta;\Delta \vdash \tau :: k$ and $\Theta;\Delta \vdash \tau :: k'$, then $k = k'$.

*Proof*
As with uniqueness of sorting, no non-variable type is compatible with multiple kinding rules. Since all kinding rules except for K-Var ascribe a specific kind, the only remaining case is for type variables. The kind environment $\Delta$ is a well-defined partial function, so $\Delta(x) = k$ and $\Delta(x) = k'$ imply $k = k'$.    □

*Lemma 2.4 (Preservation of kinds under index substitution)*
If $\Theta, x :: \gamma;\Delta \vdash \tau :: k$ and $\Theta \vdash \iota :: \gamma$ then $\Theta;\Delta \vdash \tau[x \mapsto \iota_x] :: k$.

*Proof*
This is straightforward induction on the original kind derivation.    □

*Lemma 2.5 (Preservation of kinds under type substitution)*
If $\Theta;\Delta, x :: k_x \vdash \tau :: k$ and $\Theta;\Delta \vdash \tau_x :: k_x$ then $\Theta;\Delta \vdash \tau[x \mapsto \tau_x] :: k$.

*Proof*
This is also provable by induction on the kind derivation for $\tau$.    □

### *2.3.3 Typing*

The typing rules in Figure 7 relate a full environment ($\Theta$ mapping index variables to sorts, $\Delta$ mapping type variables to kinds, and $\Gamma$ mapping term variables to types), a term (whether an atom or an expression), and its type under that environment. Since an array type might have its shape described in multiple different ways, *e.g.*, a vector of length 6 or a vector of length $1 + 5$, the T-Eqv rule makes reference to a type equivalence judgment (presented in full detail in §2.3.4) which reconciles such differences according to the algebraic theory of type indices (presented earlier in §2.2).

$$\boxed{\Theta;\Delta;\Gamma \vdash t : \tau}$$

$$\frac{}{\Theta;\Delta;\Gamma \vdash \mathfrak{o} : \mathscr{S}[\![\mathfrak{o}]\!]} \ \text{T-OP} \qquad \frac{(x : \tau) \in \Gamma}{\Theta;\Delta;\Gamma \vdash x : \tau} \ \text{T-VAR} \qquad \frac{\Theta;\Delta;\Gamma \vdash t : \tau' \qquad \tau \cong \tau'}{\Theta;\Delta;\Gamma \vdash t : \tau} \ \text{T-EQV}$$

$$\frac{\begin{array}{c} \Theta;\Delta;\Gamma \vdash \mathfrak{a}_j : \tau \ \text{for each } j \\ \Theta;\Delta \vdash \tau :: \mathtt{Atom} \\ \mathit{Length}[\![\mathfrak{a}\ldots]\!] = \prod n \ldots \end{array}}{\begin{array}{c} \Theta;\Delta;\Gamma \vdash (\mathtt{array}\ (n \ldots)\ \mathfrak{a} \ldots) \\ : (\mathtt{Arr}\ \tau\ (\mathtt{Shp}\ n \ldots)) \end{array}} \ \text{T-ARRAY} \qquad \frac{\begin{array}{c} \Theta;\Delta \vdash \tau :: \mathtt{Atom} \\ 0 \in n \ldots \end{array}}{\begin{array}{c} \Theta;\Delta;\Gamma \vdash (\mathtt{array}\ (n \ldots)\ \tau) \\ : (\mathtt{Arr}\ \tau\ (\mathtt{Shp}\ n \ldots)) \end{array}} \ \text{T-0A}$$

$$\frac{\begin{array}{c} \Theta;\Delta;\Gamma \vdash e_j : (\mathtt{Arr}\ \tau\ \iota) \ \text{for each } j \\ \Theta;\Delta \vdash (\mathtt{Arr}\ \tau\ \iota) :: \mathtt{Array} \\ \mathit{Length}[\![e\ldots]\!] = \prod n \ldots \end{array}}{\begin{array}{c} \Theta;\Delta;\Gamma \vdash (\mathtt{frame}\ (n \ldots)\ e \ldots) \\ : (\mathtt{Arr}\ \tau\ (\mathtt{++}\ (\mathtt{Shp}\ n \ldots)\ \iota)) \end{array}} \ \text{T-FRAME} \qquad \frac{\begin{array}{c} \Theta;\Delta \vdash \tau :: \mathtt{Atom} \\ \Theta \vdash \iota :: \mathtt{Shape} \qquad 0 \in n \ldots \end{array}}{\begin{array}{c} \Theta;\Delta;\Gamma \vdash (\mathtt{frame}\ (n \ldots)\ (\mathtt{Arr}\ \tau\ \iota)) \\ : (\mathtt{Arr}\ \tau\ (\mathtt{++}\ (\mathtt{Shp}\ n \ldots)\ \iota)) \end{array}} \ \text{T-0F}$$

$$\frac{\begin{array}{c} \Theta;\Delta;\Gamma,x : \tau \ldots \vdash e : \tau' \\ \Theta;\Delta \vdash \tau :: \mathtt{Array} \ \text{for each } j \end{array}}{\begin{array}{c} \Theta;\Delta;\Gamma \vdash (\lambda\ ((x\ \tau)\ \ldots)\ e) \\ : (\mathtt{->}\ (\tau \ldots)\ \tau') \end{array}} \ \text{T-LAM} \qquad \frac{\Theta;\Delta,x :: k \ldots;\Gamma \vdash v : \tau}{\begin{array}{c} \Theta;\Delta;\Gamma \vdash (\mathtt{T}\lambda\ ((x\ k)\ \ldots)\ v) \\ : (\mathtt{Forall}\ ((x\ k)\ \ldots)\ \tau) \end{array}} \ \text{T-TLAM}$$

$$\frac{\Theta,x :: \gamma \ldots;\Delta;\Gamma \vdash v : \tau}{\begin{array}{c} \Theta;\Delta;\Gamma \vdash (\mathtt{I}\lambda\ ((x\ \gamma)\ \ldots)\ v) \\ : (\mathtt{Pi}\ ((x\ \gamma)\ \ldots)\ \tau) \end{array}} \ \text{T-ILAM} \qquad \frac{\begin{array}{c} \Theta \vdash \iota :: \gamma \ \text{for each } j \\ \Theta;\Delta \vdash (\mathtt{Sigma}\ ((x\ \gamma)\ \ldots)\ \tau) :: \mathtt{Atom} \\ \Theta;\Delta;\Gamma \vdash e : \tau[x \mapsto \iota, \ldots] \end{array}}{\begin{array}{c} \Theta;\Delta;\Gamma \vdash (\mathtt{box}\ \iota \ldots e\ (\mathtt{Sigma}\ ((x\ \gamma)\ \ldots)\ \tau)) \\ : (\mathtt{Sigma}\ ((x\ \gamma)\ \ldots)\ \tau) \end{array}} \ \text{T-BOX}$$

$$\frac{\Theta;\Delta;\Gamma \vdash e : (\mathtt{Arr}\ (\mathtt{Forall}\ ((x\ k)\ \ldots)\ (\mathtt{Arr}\ \tau_u\ \iota_u))\ \iota_f) \qquad \Theta;\Delta \vdash \tau_j :: k_j \ \text{for each } j}{\Theta;\Delta;\Gamma \vdash (\mathtt{t\text{-}app}\ e\ \tau \ldots) : (\mathtt{Arr}\ \tau_u[x \mapsto \tau, \ldots]\ (\mathtt{++}\ \iota_f\ \iota_u))} \ \text{T-TAPP}$$

$$\frac{\Theta;\Delta;\Gamma \vdash e : (\mathtt{Arr}\ (\mathtt{Pi}\ ((x\ \gamma)\ \ldots)\ (\mathtt{Arr}\ \tau_p\ \iota_p))\ \iota_f) \qquad \Theta \vdash \iota_j :: \gamma_j \ \text{for each } j}{\Theta;\Delta;\Gamma \vdash (\mathtt{i\text{-}app}\ e\ \iota \ldots) : (\mathtt{Arr}\ \tau_p[x \mapsto \iota, \ldots]\ (\mathtt{++}\ \iota_f\ \iota_p[x \mapsto \iota, \ldots]))} \ \text{T-IAPP}$$

$$\frac{\begin{array}{c} \Theta;\Delta;\Gamma \vdash e_s : (\mathtt{Arr}\ (\mathtt{Sigma}\ ((x_i'\ \gamma)\ \ldots)\ \tau_s)\ \iota_s) \\ \Theta,x_i :: \gamma \ldots;\Delta;\Gamma,x_e : \tau_s[x_i' \mapsto x_i, \ldots] \vdash e_b : (\mathtt{Arr}\ \tau_b\ \iota_b) \\ \Theta;\Delta \vdash (\mathtt{Arr}\ \tau_b\ \iota_b) :: \mathtt{Array} \end{array}}{\Theta;\Delta;\Gamma \vdash (\mathtt{unbox}\ (x_i \ldots x_e\ e_s)\ e_b) : (\mathtt{Arr}\ \tau_b\ (\mathtt{++}\ \iota_s\ \iota_b))} \ \text{T-UNBOX} \qquad \text{T-UNBOX}$$

$$\frac{\begin{array}{c} \Theta;\Delta;\Gamma \vdash e_f : (\mathtt{Arr}\ (\mathtt{->}\ ((\mathtt{Arr}\ \tau\ \iota)\ \ldots)\ (\mathtt{Arr}\ \tau'\ \iota'))\ \iota_f) \\ \Theta;\Delta;\Gamma \vdash e_a : (\mathtt{Arr}\ \tau\ (\mathtt{++}\ \iota_a\ \iota))\ \ldots \qquad \iota_p = \bigsqcup\{\iota_f\ \iota_a \ldots\} \end{array}}{\Theta;\Delta;\Gamma \vdash (e_f\ e_a \ldots) : (\mathtt{Arr}\ \tau'\ (\mathtt{++}\ \iota_p\ \iota'))} \ \text{T-APP}$$

Fig. 7. Typing rules

The signature $\mathscr{S}$, referenced in the T-OP rule, is a function mapping from primitive operators to their types. For example, $\mathscr{S}[\![+]\!]$ is (`-> ((Arr Num (Shp)) (Arr Num (Shp)))` `(Arr Num (Shp))`), meaning + is an operator which consumes two scalar numbers and produces one scalar number.

Array literals and nested frames both include a length check: the number of atoms or cells must be equal to the product of the given dimensions. In the case of empty arrays,

the length matching condition is fulfilled if and only if the array has a 0 as one of its
dimensions. Term, type, and index abstractions all bind their arguments' names in the
appropriate environment.

Typing function application starts by identifying the type of the expression in function
position. It must be an array of functions, and the array's entire shape $\iota_f$ is treated as
the function frame. The function input types, also arrays, specify the element type and
cell shape for each argument. Each cell shape $\iota$ must be a suffix of the shape of the
corresponding actual argument; the remainder $\iota_a$ is the argument's frame. The maximum of
these frames under prefix ordering (where $[2\,3] \sqsubseteq [2\,3\,2]$ but $[2\,3] \not\sqsubseteq [6\,3\,2]$) is the principal
frame $\iota_p$. That is, the function and argument arrays will all be lifted so as to have $\iota_p$ as their
frames when the program runs. Then $\iota_p$ is used as the frame around the function's output
type to give the result type for this function application.

Type and index application also require arrays in function position, but they can skip
prefix comparison as type and index arguments do not come in arrays that must be split
into frames of cells. Thus the function's frame shape $\iota_f$ passes through unaltered, and
arguments are substituted into the body type $\tau_b$ to produce the resulting array's element
type.

When constructing a box, a dependent-sum type annotation is provided. The box's
index components must match their declared sorts, and substituting them into the body
of the dependent sum type must produce a type that matches the box's array component.
Unboxing requires that $e_{box}$, the expression being destructed, be a dependent sum. The
unbox form names the sum's index and array components and adds them to the sort and
type environments when checking $e_{body}$. Although the index components are in scope
while checking the body, information hidden by the existential is not permitted to leak
out: The end result type $\tau_{body}$ must be well-formed without relying on the extended sort
environment. Unboxing a frame of boxes (scalars) produces a frame of result cells, similar
to lifting function application.

Anticipating a progress lemma, we prove a canonical-forms lemma for Remora's typing
rules. Following the atom/array distinction, we have separate lemmas for atoms and arrays.
Although an atom can contain an array if that atom is a box, we avoid mutual dependence
between the lemmas by not making any claim about the syntactic structure of the box's
contents.

*Lemma 2.6 (Canonical forms for atomic values)*
Let $\mathfrak{v}$ be a well-typed atomic value, that is, $\cdot;\cdot;\cdot \vdash \mathfrak{v} : \tau$.

1. If $\tau$ is of the form $(\text{->} (\tau_i \dots) \tau_o)$,
   then $\mathfrak{v}$ is of the form $\mathfrak{o}$ or $(\lambda ((x\,\tau_i) \dots) e)$.
2. If $\tau$ is of the form $(\text{Forall} ((x\,k) \dots) \tau_u)$,
   then $\mathfrak{v}$ is of the form $(\text{T}\lambda ((x_u\,k) \dots) v)$.
3. If $\tau$ is of the form $(\text{Pi} ((x\,\gamma) \dots) \tau_p)$,
   then $\mathfrak{v}$ is of the form $(\text{I}\lambda ((x_p\,\gamma) \dots) v)$.
4. If $\tau$ is of the form $(\text{Sigma} ((x\,\gamma) \dots) \tau_b)$,
   then $\mathfrak{v}$ is of the form $(\text{box}\ \iota \dots v_b\ (\text{Sigma} ((x_b\,\gamma) \dots) \tau'_b))$,
   with $\tau \cong (\text{Sigma} ((x_b\,\gamma) \dots) \tau'_b)$.

$$\boxed{\tau \cong \tau'}$$

$$\frac{}{\tau \cong \tau} \text{ TEQV-REFL} \qquad \frac{\tau \cong \tau' \qquad \text{VALID} \llbracket \iota \equiv \iota' \rrbracket}{(\text{Arr } \tau \, \iota) \cong (\text{Arr } \tau' \, \iota')} \text{ TEQV-ARRAY}$$

$$\frac{\tau_{ij} \cong \tau'_{ij} \text{ for each } j \qquad \tau_o \cong \tau'_o}{(\text{-> } (\tau_i \dots) \, \tau_o) \cong (\text{-> } (\tau'_i \dots) \, \tau'_o)} \text{ TEQV-FN}$$

$$\frac{\tau[x \mapsto x_f, \dots] \cong \tau'[x' \mapsto x_f, \dots] \quad \text{with fresh } x_f \dots}{(\text{Forall } ((x \, k) \dots) \, \tau) \cong (\text{Forall } ((x' \, k) \dots) \, \tau')} \text{ TEQV-UNIV}$$

$$\frac{\tau[x \mapsto x_f, \dots] \cong \tau'[x' \mapsto x_f, \dots] \quad \text{with fresh } x_f \dots}{(\text{Pi } ((x \, \gamma) \dots) \, \tau) \cong (\text{Pi } ((x' \, \gamma) \dots) \, \tau')} \text{ TEQV-PI}$$

$$\frac{\tau[x \mapsto x_f, \dots] \cong \tau'[x' \mapsto x_f, \dots] \quad \text{with fresh } x_f \dots}{(\text{Sigma } ((x \, \gamma) \dots) \, \tau) \cong (\text{Sigma } ((x' \, \gamma) \dots) \, \tau')} \text{ TEQV-SIGMA}$$

Fig. 8. Type equivalence

5. If $\tau$ is of the form $B$,
then $\mathfrak{v}$ is of the form $\mathfrak{b}$.

*Proof*
The type derivation may end with T-EQV, so we consider the subderivation prior to all final T-EQV instances. We then examine which typing rules can ascribe a type of the right form and the identify what form the term must take to match those rule. $\square$

*Lemma 2.7* (*Canonical forms for arrays*)
Let $v$ be a well-typed value, that is, $\cdot; \cdot; \cdot \vdash v : \tau$,

1. If $\tau$ is of the form $(\text{Arr } (\text{-> } (\tau_i \dots) \, \tau_o) \, \iota)$,
then $v$ is of the form $(\text{array } (n \dots) \, \mathfrak{f} \dots)$.
2. If $\tau$ is of the form $(\text{Arr } (\text{Forall } ((x \, k) \dots) \, \tau_u) \, \iota)$,
then $v$ is of the form $(\text{array } (n \dots) \, (\text{T}\lambda \, ((x_u \, k) \dots) \, v_u) \dots)$.
3. If $\tau$ is of the form $(\text{Arr } (\text{Pi } ((x \, \gamma) \dots) \, \tau_p) \, \iota)$,
then $v$ is of the form $(\text{array } (n \dots) \, (\text{I}\lambda \, ((x_p \, \gamma) \dots) \, v_u) \dots)$.
4. If $\tau$ is of the form $(\text{Arr } (\text{Sigma } ((x \, \gamma) \dots) \, \tau_b) \, \iota)$,
then $v$ is of the form $(\text{array } (n \dots) \, (\text{box } \iota \dots \, v_b \, (\text{Sigma } ((x_b \, \gamma) \dots) \, \tau_b)) \dots)$,
with $\tau \cong (\text{Sigma } ((x_b \, \gamma) \dots) \, \tau'_b)$.
5. If $\tau$ is of the form $(\text{Arr } B \, \iota)$,
then $v$ is of the form $(\text{array } (n \dots) \, \mathfrak{b} \dots)$,
with $\cdot; \cdot; \cdot \vdash \mathfrak{b} : B$ for each of $\mathfrak{b} \dots$.

*Proof*
This proceeds like the proof for Lemma 2.6. $\square$

### 2.3.4 Type equivalence

Remora's typing rules rely on a type-equivalence relation defined in Figure 8. The equivalence relation is essentially $\alpha$-equivalence augmented with a check as to whether array

shapes are guaranteed to be equal. We expect the relation $\cong$ actually to be an equivalence relation, *i.e.*, reflexive, symmetric, and transitive. Only reflexivity has its own inference rule, so we now show symmetry and transitivity.

*Lemma 2.8* (*Symmetry of $\cong$*)
If $\tau \cong \tau'$, then $\tau' \cong \tau$.

*Proof*
This follows via straightforward induction on the derivation of $\tau_0 \cong \tau_1$.     □

*Lemma 2.9* (*Transitivity of $\cong$*)
If $\tau_0 \cong \tau_1$ and $\tau_1 \cong \tau_2$, then $\tau_0 \cong \tau_2$.

*Proof*
This follows from induction on the derivations of $\tau_0 \cong \tau_1$ and $\tau_1 \cong \tau_2$. Since both derivations $\tau_1$, the structure of the equivalence rules prohibits the derivations from ending with different rules (other than TEQV-REFL, which passes that structural requirement on to its premises).     □

*Theorem 2.1*
$\cong$ is an equivalence relation.

A type-equivalence relation should not cross kind boundaries. Violation of this principle would allow use of T-EQV to ascribe an ill-kinded type to a well-typed term. It follows directly from inspection of the equivalence rules that they will not relate an `Atom` with an `Array`, but correct use of type and index variables remains to be proven. To that end, we show that two equivalent types will be ascribed the same kind by the same environment.

*Lemma 2.10*
If $\Theta;\Delta \vdash \tau :: k$ and $\tau \cong \tau'$, then $\Theta;\Delta \vdash \tau' :: k$.

*Proof*
This result is proven induction on the derivation of $\tau \cong \tau'$. In each case, the induction hypothesis converts a kind derivation for some fragment of $\tau$ into a kind derivation for a corresponding fragment of $\tau'$ (and similar for index fragments).     □

We also expect type equivalence to be well-behaved under substitution. Ultimately, substituting equivalent types or indices into equivalent types ought to produce equivalent types. Proving that result by induction on derivation of equivalence is straightforward except for the REFL case.

*Lemma 2.11*
If VALID $[\![ \iota \equiv \iota' ]\!]$, then for any index variable $x$, $\tau[x \mapsto \iota] \cong \tau[x \mapsto \iota']$.

*Proof*
This is provable using induction on the structure of $\tau$. Only the case for arrays makes direct use of $\iota$ and $\iota'$; the other cases simply use the induction hypothesis to prove the premises of the derivation of $\tau[x \mapsto \iota] \cong \tau[x \mapsto \iota']$.     □

*Lemma 2.12*
If $\tau_x \cong \tau'_x$, then for any type variable $x$, $\tau[x \mapsto \tau_x] \cong \tau[x \mapsto \tau'_x]$.

*Proof*

We use induction on the structure of $\tau$. The cases for universals, dependent products, and dependent sums require instantiating the induction hypothesis with a substitution of fresh type variables $x_f \ldots$. For example, when $\tau = (\texttt{Forall } ((x_u\ k)\ \ldots )\ \tau_u)$, the induction hypothesis promises the equivalence of $\tau_u$ after substituting in $x_f \ldots$ for $x_u \ldots$ and *also* $\tau_x$ or $\tau'_x$ for $x$. $\qquad\square$

*Theorem 2.2*

If $\tau \cong \tau'$ and $\tau_x \cong \tau'_x$, then for any type variable $x$, $\tau[x \mapsto \tau_x] \cong \tau'[x \mapsto \tau'_x]$.

*Proof*

We use induction on the derivation of $\tau \cong \tau'$. In each case, the induction hypothesis provides equivalence derivations for corresponding fragments of $\tau[x \mapsto \tau_x]$ and $\tau'[x \mapsto \tau'_x]$, which can then be used to prove the substituted types themselves equivalent. $\qquad\square$

Having defined the typing judgment and the type-equivalence relation on which it builds, we can now prove the usual results about typing in Remora. The T-EQV rule can allow many types to be ascribed to a single term, but we will prove that an environment and term can only map to a single equivalence class.

*Theorem 2.3* (*Uniqueness of typing, up to equivalence*)

If $\Theta;\Delta;\Gamma \vdash t : \tau$ and $\Theta;\Delta;\Gamma \vdash t : \tau'$, then $\tau \cong \tau'$.

*Proof*

This can be proven by induction on $t$, showing that all derivations of $\Theta;\Delta;\Gamma \vdash t : \tau'$ must end with the same non-T-EQV rule (chosen according to the structure of $t$) followed by 0 or more T-EQV instances, which keeps the result in the same equivalence class as $\tau$. $\qquad\square$

We also require guarantees about substitution in terms: replacing an index variable with an appropriately-sorted index, a type variable with an appropriately-kinded type, or a term variable with an appropriately-typed expression should not change the type of the original term. If substitution turns a term $t$ with type $\tau$ into $t'$ with type $\tau'$, where $\tau \cong \tau'$, we can add a T-EQV at the end of the new type derivation to conclude $t'$ has type $\tau$. As such, we do not need to include an "up to equivalence" caveat when stating the preservation of typing lemmas.

*Lemma 2.13* (*Preservation of types under index substitution*)

If $\Theta, x :: \gamma; \Delta; \Gamma \vdash t : \tau$ and $\Theta \vdash \iota :: \gamma$ then $\Theta; \Delta; \Gamma[x \mapsto \iota_x] \vdash t[x \mapsto \iota_x] : \tau[x \mapsto \iota_x]$.

*Proof*

This is straightforward induction on the derivation of $\Theta, x :: \gamma; \Delta; \Gamma \vdash t : \tau$. $\qquad\square$

*Lemma 2.14* (*Preservation of types under type substitution*)

If $\Theta; \Delta, x :: k; \Gamma \vdash t : \tau$ and $\Theta; \Delta \vdash \tau_x :: k$ then $\Theta; \Delta; \Gamma[x \mapsto \tau_x] \vdash t[x \mapsto \tau_x] : \tau[x \mapsto \tau_x]$.

*Proof*

This is straightforward induction on the derivation of $\Theta; \Delta, x :: k; \Gamma \vdash t : \tau$. $\qquad\square$

*Lemma 2.15* (*Preservation of types under term substitution*)

If $\Theta; \Delta; \Gamma, x : \tau_x \vdash t : \tau$ and $\Theta; \Delta; \Gamma \vdash e_x : \tau_x$ then $\Theta; \Delta; \Gamma \vdash t[x \mapsto e_x] : \tau$.

*Proof*
We use induction on the derivation of $\Theta;\Delta;\Gamma, x : \tau_x \vdash t : \tau$. $\qquad \Box$

We call an environment well-formed, noted as $\Theta;\Delta \vdash \Gamma$, if for every binding $x : \tau \in \Gamma$, we can derive $\Theta;\Delta \vdash \tau :: \texttt{Array}$. This is the expected case, rather than permitting $\tau$ to have kind $\texttt{Atom}$, because a lone variable is an expression and ought to stand for an array value.

When we show that the typing judgment only ascribes types of the appropriate kind, the case for the T-EQV rule relies on the earlier lemma that the type equivalence relation respects kinding, *i.e.*, two equivalent types will have the same kind when checked in the same environment.

*Theorem 2.4* (*Ascription of well-kinded types*)
Given $\Theta;\Delta;\Gamma \vdash t : \tau$ where $\Theta;\Delta \vdash \Gamma$:

- If $t$ is an expression, then $\Theta;\Delta \vdash \tau :: \texttt{Array}$
- If $t$ is an atom, then $\Theta;\Delta \vdash \tau :: \texttt{Atom}$

*Proof*
This follows from induction on the derivation of $\Theta;\Delta;\Gamma \vdash t : \tau$. It is not sufficient to point out that each typing rule ascribes a type whose form matches the appropriate kind. Elimination-form cases call for a little extra work. For UNBOX, the kind check on the result type is necessary to ensure that existentially quantified variables do not leak out. The APP case must ensure that index variables in the ascribed type actually appear in the environment. This is guaranteed because the principal frame is always chosen to be one of the function- or argument-position frames. $\qquad \Box$

### *2.4 Dynamic Semantics*

In the dynamic semantics for Remora, the way function application is lifted to work on aggregate data depends on the types of the function and argument terms. Consulting type information avoids a "hole" in the semantics of untyped array-oriented code, where a frame whose shape includes a 0 dimension evaluates to an array with indeterminate shape— there are no concrete cells whose shape can be used to determine the overall shape of the resulting array. Instead, the function's type tells us the shape of the resulting cells, even when there are zero such cells.

The small-step operational semantics, given in Figure 10, assumes every atom or expression has been tagged with its type. For example, $\beta$-reduction requires that each atom in the function position array have input types $\tau \ldots$ and that the argument arrays' types also match $\tau \ldots$. This matching is still subject to the type equivalence rules described in §2.3, *e.g.*, a function tagged as having input type (Array Int (++ (Shp 3) (Shp 4))) can be applied to an argument tagged with type (Array Int (Shp 3 4)). Because every term now has type annotations attached, we drop the "empty" array and frame syntactic forms. Their replacements use the standard array and frame syntax with an empty list of atoms or cells, and the atom or cell type is implied by the expression's type annotation.

Several list-processing metafunctions are used in defining the reduction rules. These metafunctions are defined in Figure 9. *Split$_n$* turns a list into a list of lists, made up of the consecutive length-$n$ pieces of the original list. For example, *Split$_3$* $[\![(1\ 2\ 3\ 4\ 5\ 6)]\!]$ is

$$Split_n \, [\![(a_1,\ldots,a_m)]\!] = ((a_1,\ldots,a_n),(a_{n+1},\ldots,a_{2n}),\ldots,(a_{m-n+1},\ldots,a_m))$$

$$Rep_n \, [\![(a_1,\ldots,a_m)]\!] = (a_{1,1},\ldots,a_{1,n},\ldots,a_{m,1},\ldots,a_{m,n}) \quad \text{where } a_{i,j} = a_i$$

$$Concat \, [\![((a_{1,1},\ldots,a_{1,n}),\ldots,(a_{m,1},\ldots,a_{m,n}))]\!] = (a_{1,1},\ldots,a_{1,n},\ldots,a_{m,1},\ldots,a_{m,n})$$

$$Transpose((a_{1,1},\ldots,a_{1,n}),\ldots,(a_{m,1},\ldots,a_{m,n})) = ((a_{1,1},\ldots,a_{m,1}),\ldots,(a_{n,1},\ldots,a_{m,n}))$$

Fig. 9. List-processing metafunctions

((1 2 3) (4 5 6)). *Concat* flattens a list of lists into a single list, effectively reversing a *Split*. $Rep_n$ constructs a new list by repeating each element of the original list $n$ times. $Rep_2 \, [\![(0 \; 1)]\!]$ is (0 0 1 1). Used on nested lists, the inner lists are treated atomically: $Rep_2 \, [\![((1 \; 2 \; 3) \; (4 \; 5 \; 6))]\!]$ is ((1 2 3) (1 2 3) (4 5 6) (4 5 6)). *Transpose* takes a list of lists, where the inner lists all have the same length, and produces a new list of lists whose $i^{th}$ element contains the $i^{th}$ elements of each original inner list.

The reduction rules themselves are given in Figure 10. Remora's function application is split into stages for replicating cells to make frame shapes match (*lift*), mapping the functions to corresponding argument cells (*map*), and gathering the result cells back into an array (*collapse*).

Performing a *lift* step identifies the function array's frame, the sequence $[n_f \ldots]$, and each argument's frame, $[n_a \ldots]$. Then the sequence $[n_p \ldots]$ is chosen to be the largest frame according to prefix ordering. We require that at least one function or argument frame be different from the principal frame—otherwise, a *map* step would be appropriate instead. Each argument's cell size $n_{ac}$ is the product of the dimensions $[n_{in} \ldots]$ of the function's input type at that position; the function array's cell size is always 1. The number of replicas needed for each cell ($n_{fe}$ for the function and $n_{fa}$ for each argument) is determined by multiplying the dimensions that must be added to each corresponding frame to produce the principal frame, *i.e.*, the principal frame minus whatever prefix was already present in the original array's shape. Given these numbers, we split each array's atom list into its cells, replicate those cells to match the new array shapes, and then concatenate each array's replicated cells to produce the new function and argument arrays. Type annotations on the individual arrays update to reflect their new shapes, but the application form's type remains unchanged.

A *map* step is possible when every piece of a function application has the same frame shape. Then the application becomes a `frame` of application forms, which themselves all have scalar principal frame. This requires breaking each argument array's atom list into its individual cells' atom lists, then transposing to match the first cell of each argument with the first function, the second cell of each argument with the second function, and so on.

When function application has a scalar in function position, and every argument array matches the function's corresponding input type, then we can $\beta$-reduce or $\delta$-reduce. $\beta$-reduction performs conventional $\lambda$-calculus substitution. The $\delta$ rule uses a family of metafunctions, each associated with a primitive operator. No `frame` construct is necessary in either result, as this is the degenerate case of function lifting—the principal frame is scalar.

Applying type and index abstractions is handled by the $t\beta$ and $i\beta$ rules. The application frame is the shape of the array of type or index abstractions, since there are no argument

*arrays*. Every $\text{T}\lambda$ or $\text{I}\lambda$ is applied to the full list of type or index arguments. Substitution into the body of each abstraction should be read as affecting type annotations as well as subterms: if we are replacing the type variable T with Int, then $\text{x}^{(\text{Array T (Shp 3)})}$ becomes $\text{x}^{(\text{Array Int (Shp 3)})}$.

Once a `frame` has every one of its cells reduced to an `array` literal, the nested representation can be merged into a single literal. In the case where one of $n \ldots$ is 0, there will be no cells to examine to determine the cell dimensions $n' \ldots$, so this information is taken from the type annotation on the `frame` form. The type annotation itself passes through unchanged. The atom lists from the cells are concatenated to produce the collapased array's atom list.

Destructing a `box` with an `unbox` form behaves like a conventional `let`. The result is the body $e$, where the index variables $x_i \ldots$ are replaced with the box's indices $\iota \ldots$, and the term variable $x_e$ is replaced with the contained array $v$.

### *2.5 Type Soundness*

The value of a type soundness theorem for Remora is not only assurance that well-typed programs do not suffer from shape-mismatching errors. It also ensures that the types ascribed to program terms accurately describe the shapes of the data those terms compute. That is the guarantee that justifies a compiler's use of the type system as a static analysis for array shape.

With supporting lemmas, such as canonical forms and substitution, already taken care of, we now establish progress and preservation lemmas. Since we have not committed to a collection of primitive operators that are all total functions, the progress lemma acknowledges the possibility of non-shape errors, such as division by zero. However, we do assume that any value returned by a primitive operator inhabits that operator's output type.

*Lemma 2.16* (*Progress*)
Given an expression $e$ such that $\cdot; \cdot; \cdot \vdash e : \tau$, one of the following holds:

- $e$ is a value $v$
- There exists $e'$ such that $e \mapsto e'$
- $e$ is $\mathbb{E}\big[((\text{array ()}\ \mathfrak{o})\ v \ldots)\big]$ where $\mathfrak{o}$ is a partial function applied to appropriately-typed values outside its domain.

*Proof*
We use induction on the derivation of $\cdot; \cdot; \cdot \vdash e : \tau$. We consider only cases for typing rules which apply to expressions (as opposed to atoms). Since we do not reduce under a binder, our assumed type derivation ensures that the reducible subexpression of $e$ is also typable using an empty environment.

An `array` form which is not already a value must have some non-value atom. That atom must itself contain a non-value expression, with its own type derivation. So the induction hypothesis implies that it can take a reduction step or is a mis-applied primitive operator. Similar reasoning applies to `frame` forms: either we have a *collapse* redex, or some cell subexpression in the frame can make progress.

$$((\texttt{array}\ (n_f \ldots)\ \mathfrak{v}_f \ldots)^{(\texttt{Arr}\ (\texttt{->}\ ((\texttt{Arr}\ \tau_i\ (\texttt{Shp}\ n_i \ldots))\ldots)\ \tau_o)\ (\texttt{Shp}\ n_f \ldots))}$$
$$(\texttt{array}\ (n_a \ldots n_i \ldots)\ \mathfrak{v}_a \ldots)^{(\texttt{Arr}\ \tau_i\ (\texttt{Shp}\ n_a \ldots n_i \ldots))}\ \ldots)$$
$$\mapsto_{lift}$$
$$((\texttt{array}\ (n_p \ldots)$$
$$Concat\ \Big[\!\!\Big[ Rep_{n_{fe}}\ \big[\!\!\big[ Split_1\ [\!\![ \mathfrak{v}_f \ldots ]\!\!] \big]\!\!\big] \Big]\!\!\Big])^{(\texttt{Arr}\ (\texttt{->}\ ((\texttt{Arr}\ \tau_i\ (\texttt{Shp}\ n_i \ldots))\ldots)\ \tau_o)\ (\texttt{Shp}\ n_p \ldots))}$$
$$(\texttt{array}\ (n_p \ldots n_i \ldots)$$
$$Concat\ \Big[\!\!\Big[ Rep_{n_{ae}}\ \big[\!\!\big[ Split_{n_{ac}}\ [\!\![ \mathfrak{v}_a \ldots ]\!\!] \big]\!\!\big] \Big]\!\!\Big])^{(\texttt{Arr}\ \tau_i\ (\texttt{Shp}\ n_p \ldots n_i \ldots))}\ \ldots)$$

where

Not all of $(n_f \ldots), (n_a \ldots) \ldots$ are equal

$$n_p \ldots \ =\ \bigsqcup [\!\![ (n_f \ldots)\ (n_a \ldots)\ldots ]\!\!] \qquad n_{fe}\ =\ \prod \frac{(n_p \ldots)}{(n_f \ldots)}$$
$$n_{ae} \ldots \ =\ \prod \frac{(n_p \ldots)}{(n_a \ldots)} \ldots \qquad\quad n_{ac} \ldots \ =\ \left(\prod (n_i \ldots)\right) \ldots$$

$$((\texttt{array}\ (n_f \ldots)\ \mathfrak{v}_f \ldots)^{(\texttt{Arr}\ (\texttt{->}\ ((\texttt{Arr}\ \tau_i\ (\texttt{Shp}\ n_i \ldots))\ldots)\ \tau_o)\ (\texttt{Shp}\ n_f \ldots))}$$
$$(\texttt{array}\ (n_f \ldots\ n_i \ldots)\ \mathfrak{v}_a \ldots)^{(\texttt{Arr}\ \tau_i\ (\texttt{Shp}\ n_f \ldots\ n_i \ldots))}\ \ldots)$$
$$\mapsto_{map}$$
$$(\texttt{frame}\ (n_f \ldots)$$
$$((\texttt{array}\ ()\ \mathfrak{v}_f)^{(\texttt{Arr}\ (\texttt{->}\ ((\texttt{Arr}\ \tau_i\ (\texttt{Shp}\ n_i \ldots))\ldots)\ \tau_o)\ (\texttt{Shp}))}$$
$$(\texttt{array}\ (n_i \ldots)\ \mathfrak{v}_c \ldots)^{(\texttt{Arr}\ \tau_i\ (\texttt{Shp}\ n_i \ldots))}\ \ldots)^{\tau_o}\ \ldots)$$

where
$$n_c \ldots \ =\ \left(\prod n_i \ldots\right) \ldots$$
$$((\mathfrak{v}_c \ldots) \ldots) \ldots \ =\ Transpose\ \big[\!\!\big[ Split_{n_c}\ [\!\![ \mathfrak{v}_a \ldots ]\!\!] \ldots \big]\!\!\big]$$
$$Length\ [\!\![ n_f \ldots ]\!\!]\ >\ 0$$

$$((\texttt{array}\ ()\ (\lambda\ ((x\ \tau)\ \ldots)\ e))\ v^\tau \ldots)$$
$$\mapsto_\beta\ e[x \mapsto v^\tau, \ldots]$$

$$(\texttt{t-app}\ (\texttt{array}\ (n \ldots)\ (\texttt{T}\lambda\ ((x\ k)\ \ldots)\ e)\ \ldots)\ \tau \ldots)$$
$$\mapsto_{t\beta}\ (\texttt{frame}\ (n \ldots)\ e[x \mapsto \tau, \ldots]\ \ldots)$$

$$(\texttt{i-app}\ (\texttt{array}\ (n \ldots)\ (\texttt{I}\lambda\ ((x\ \gamma)\ \ldots)\ e)\ \ldots)\ \iota \ldots)$$
$$\mapsto_{i\beta}\ (\texttt{frame}\ (n \ldots)\ e[x \mapsto \iota, \ldots]\ \ldots)$$

$$(\texttt{frame}\ (n \ldots)\ (\texttt{array}\ (n' \ldots)\ \mathfrak{v}\ \ldots)\ \ldots)^{(\texttt{Arr}\ \tau\ (\texttt{Shp}\ n \ldots n' \ldots))}$$
$$\mapsto_{collapse}\ (\texttt{array}\ (n \ldots n' \ldots)\ Concat\ [\!\![ (\mathfrak{v} \ldots) \ldots ]\!\!])^{(\texttt{Arr}\ \tau\ (\texttt{Shp}\ n \ldots n' \ldots))}$$

$$(\texttt{unbox}\ (x_i \ldots\ x_e\ (\texttt{array}\ (n_s \ldots)\ (\texttt{box}\ \iota \ldots\ v\ \tau) \ldots))\ e)$$
$$\mapsto_{unbox}\ (\texttt{frame}\ (n_s \ldots)\ e[x_i \mapsto \iota, \ldots, x_e \mapsto v])$$

Fig. 10. Dynamic semantics for Remora

An unbox form can either make progress in the box position (via the induction hypothesis) or take an *unbox* step. Similarly, a type or index applications can make progress in function position or take $t\beta$ or $i\beta$ step.

The function application case splits into subcases depending on whether the function and argument arrays are fully reduced and if so what their frame shapes are. If they are all

value forms, we have all scalar frames (a $\beta$ or $\delta$ redex) or all identical non-scalar frames (a *map* redex), or non-identical prefix-compatible frames (a *lift* redex). Prefix-incompatible frames are ruled out by the type derivation.  □

*Lemma 2.17* (*Preservation*)
Let $\Theta, \Delta, \Gamma$ be a well-formed environment, *i.e.*, $\Theta; \Delta \vdash \Gamma$. If $\Theta; \Delta; \Gamma \vdash e : \tau$ and $e \mapsto e'$ then $\Theta; \Delta; \Gamma \vdash e' : \tau$.

*Proof*
We use induction on the derivation of $\cdot; \cdot; \cdot \vdash e : \tau$. An `array` form which can take a reduction step must contain a reducible subexpression. Many typing rules give rise to subcases where the $e$ itself is not a redex but contains some subexpression $e_r$ which steps to $e_r'$. In these situations, the typing derivation for $e_r$ is included in that for $e$, so replacing that subderivation with one for $e_r'$ (deriving the same type, according to the induction hypothesis) produces a derivation of $\cdot; \cdot; \cdot \vdash e' : \tau$.

The remaining nontrivial subcases each correspond to particular reduction rules. As in proving Progress, the T-APP case is split into subcases based on the function and argument frames. When frames are non-identical but prefix-compatible, the resulting *lift* reduction produces an application form with the same principal frame and thus the same result type. When we have identical non-scalar frames, the *map* reduction produces a `frame` form whose frame shape is equal to the application form's principal frame and whose cell shape and atom type is the same as the function's return shape and atom type. This gives it a type equivalent to that of the *map* redex. With a scalar principal frame, we have a $\delta$ redex (trivial) or $\beta$ redex (follows from Lemma 2.15, preservation of types under term substitution). Reasoning for type and index application forms is similar (via Lemma 2.13 and Lemma 2.14 respectively). A reducible `unbox` form also substitutes a value in for a variable which is intended to have the same type, so Lemma 2.15 again ensures the result type is $\tau$.  □

*Theorem 2.5* (*Type soundness*)
If $\cdot; \cdot; \cdot \vdash e : \tau$, then either $e$ diverges, there exists $v$ such that $e \mapsto^* v$ and $\cdot; \cdot; \cdot \vdash v : \tau$, or there exist partial function $\mathfrak{o}$ and appropriately-typed arguments $v \dots$ such that $e \mapsto^* \mathbb{E}[((\texttt{array ()}\ \mathfrak{o})\ v \dots )]$.

*Proof*
We argue coinductively using the sequence of reduction steps from $e$. For any well-typed $e$, Progress (Lemma 2.16) implies that either $e$ has the form $v$, $e$ has the form $\mathbb{E}[((\texttt{array} ()\ \mathfrak{o})\ v\ \dots )]$, or $e \mapsto e'$. In the first case, the reduction sequence terminates in a value, so we have $e \mapsto^* v$. Furthermore, Preservation (Lemma 2.17) implies that $\cdot; \cdot; \cdot \vdash v : \tau$ In the second case, the reduction sequence terminates in a mis-applied primitive operator. In the third case, Preservation implies that $\cdot; \cdot; \cdot \vdash e' : \tau$.  □

# 3 Type Erasure

The dynamic semantics given in Section 2 relies on ubiquitous type annotations in order to determine how function application will proceed or how a frame of sub-arrays should collapse to a single array. While the possible case of constructing a frame with no actual

result cells whose shape can be inspected can only be resolved by consulting a type annotation, the types themselves contain more information than is strictly needed. For example, it does not matter whether we are collapsing an empty frame of functions, an empty frame of integers, or an empty frame of boxes. The result shape is the same, regardless of the type of the atoms contained within the cells. All we truly need is the resulting shape (alternatively, the result cells' shape). Similarly, evaluating a function application requires knowing the expected cell shapes for the arguments, but it could, in principle, be done without knowing anything about their atoms. Function application is still tagged with a result shape, again to head off issues arising from mapping over an empty frame.

In a type-erased version of Remora, we only need the term and index levels—the syntactic class of types is discarded. The syntax for erased Remora is given in Figure 11. Note that the grammar of type indices from Figure 1 is still in use here, although expressions, atoms, and their corresponding function and value-form subsets are now replaced with type-erased versions.

Evaluation in erased Remora proceeds similarly to explicit Remora. A function-application form has a principal frame chosen to be the largest of the function and argument frames, and a *lift* reduction replicates the function and argument arrays' atoms to bring all of the frames into agreement. The argument frames themselves are identified based on the individual argument positions' cell-shape annotations, rather than by inspecting a type annotation on the array in function position. A *map* reduction turns an application form where all pieces have the same frame into a `frame` form, where the end-result shape matches the result shape tag on the original application. Index application also maps over an array of index functions, producing a `frame` of substituted function bodies. Since the type level has been eliminated, there are no $t\lambda$ and `t-app` forms and no need for a $t\beta$ reduction rule.

The translation from explicit Remora to erased Remora consists of three erasure functions: $\mathscr{E}[\![\cdot]\!] : Expr \to \widehat{Expr}$, $\mathscr{A}[\![\cdot]\!] : Atom \to \widehat{Atom}$, and $\mathscr{T}[\![\cdot]\!] : Type \to Index$. These functions are defined in Figure 13.

We also define $\mathscr{C}[\![\cdot]\!] : Ctxt \to \widehat{Ctxt}$, given in Figure 14, which is not needed for defining the erased form of an explicit Remora program but is useful for demonstrating their equivalence.

Types in explicit Remora are turned into indices in erased Remora. These indices are the dynamic residue of types, in the same sense that term-level values are dynamic, though the are still subject to a static discipline which governs their values and their relation to the array values they describe. Array types become just the shapes used to construct them, whereas functions, universals, dependent sums and products, and base types become the "scalar" shape. Extracting the index components of all types means that type variables can be turned into index variables, which will stand for the index component of whatever type the variable originally stood for. This translation captures exactly the information that a `frame` form needs in the event that there are no cells. By extension, the term and index application forms also get the bookkeeping information needed by the `frames` they will eventually become.

For example, consider a function term whose type is `(-> (s (Arr t (Shp))) (Arr t (Shp k)))`, where `s`, `t`, and `k` are bound as `Array`, `Atom`, and `Dim` respectively. This function produces a vector of some statically uncertain length containing atoms of uncer-

$$
\begin{array}{llr}
\widehat{e} \in \widehat{Expr} & ::= & \textit{Type-erased expressions} \\
& x & \textit{Variable reference} \\
& \mid\ (\texttt{array}\ (n \ldots)\ \widehat{\mathfrak{a}} \ldots) & \textit{Array, containing atoms} \\
& \mid\ (\texttt{frame}\ \iota\ \widehat{e} \ldots) & \textit{Frame, containing sub-arrays} \\
& \mid\ (\widehat{e}_f\ (\widehat{e}_a\ \iota_a) \ldots \iota_r) & \textit{Term application} \\
& \mid\ (\texttt{i-app}\ \widehat{e}_f\ \iota_a \ldots \iota_r) & \textit{Index application} \\
& \mid\ (\texttt{unbox}\ (x_i \ldots x_e\ \widehat{e}_s)\ \widehat{e}_b\ \iota_b) & \textit{Let-binding box contents} \\
\widehat{\mathfrak{a}} \in \widehat{Atom} & ::= & \textit{Type-erased atoms} \\
& \mathfrak{b} & \textit{Base value} \\
& \mid\ \widehat{\mathfrak{f}} & \textit{Function} \\
& \mid\ (\texttt{I}\lambda\ (x \ldots)\ \widehat{v}) & \textit{Index abstraction} \\
& \mid\ (\texttt{box}\ \iota \ldots \widehat{e}) & \textit{Boxed array} \\
\widehat{\mathfrak{f}} \in \widehat{Func} & ::= & \textit{Type-erased functions} \\
& \mid\ \mathfrak{o} & \textit{Primitive operator} \\
& \mid\ (\lambda\ (x \ldots)\ \widehat{e}) & \textit{Term abstraction} \\
\widehat{v} \in \widehat{Val} & ::= & \textit{Type-erased values} \\
& x & \\
& \mid\ (\texttt{array}\ (n \ldots)\ \widehat{\mathfrak{v}} \ldots) & \\
\widehat{\mathfrak{v}} \in \widehat{Atval} & ::= & \textit{Type-erased atomic values} \\
& \mathfrak{b} & \\
& \mid\ \widehat{\mathfrak{f}} & \\
& \mid\ (\texttt{I}\lambda\ (x \ldots)\ \widehat{v}) & \\
& \mid\ (\texttt{box}\ \iota \ldots \widehat{v}) & \\
\widehat{\mathbb{E}} \in \widehat{Ctxt} & ::= & \textit{Type-erased evaluation contexts} \\
& \square & \\
& \mid\ (\texttt{array}\ (n \ldots)\ \widehat{\mathfrak{v}} \ldots (\texttt{box}\ \iota \ldots \widehat{\mathbb{E}})\ \widehat{\mathfrak{a}} \ldots) & \\
& \mid\ (\texttt{frame}\ \iota\ \widehat{v} \ldots \widehat{\mathbb{E}}\ \widehat{e} \ldots) & \\
& \mid\ (\widehat{\mathbb{E}}\ (\widehat{e}_a\ \iota_a) \ldots \iota_r) & \\
& \mid\ (\widehat{e}_f\ (\widehat{v}_a\ \iota_a) \ldots (\widehat{\mathbb{E}}\ \iota_a)\ (\widehat{e}_a\ \iota_a) \ldots \iota_r) & \\
& \mid\ (\texttt{i-app}\ \widehat{\mathbb{E}}\ \iota_a \ldots \iota_r) & \\
& \mid\ (\texttt{unbox}\ (x_i \ldots x_e\ \widehat{\mathbb{E}})\ \widehat{e}_b\ \iota_b) & \\
& \mid\ (\texttt{unbox}\ (x_i \ldots x_e\ \widehat{v}_s)\ \widehat{\mathbb{E}}\ \iota_b) & \\
\widehat{t} \in \widehat{Term} & ::= \widehat{e} \mid \widehat{\mathfrak{a}} & \textit{Type-erased terms}
\end{array}
$$

Fig. 11. Abstract syntax for type-erased Remora

tain type. When we apply this function, the explicitly typed application form describes the resulting array's type. If our arguments are a single s and a $n \times 4$ matrix of numbers, with n also bound as a Dim, the principal frame shape is (Shp n 4). So we will have result type (Arr Num (Shp n 4 k)). Type-erasing the application form must still preserve enough information to produce an array of the correct shape, even if n turns out to be 0, leaving us with no result cells whose shape we can inspect. However, the dynamic semantics does not rely on knowing that the result array contains Nums. The binders for index variables n and k, which must be either I$\lambda$ or unbox, are still present in the type-erased program, since the indices they eventually bind to those variables will affect the program's semantics. The T$\lambda$s which bind s and t turn into I$\lambda$s, though the variable t is never used in the type-erased program. If any type argument was bound to s in the original program, we replace it with

$$((\text{array } (n_f \dots) \; \widehat{\mathfrak{v}}_f \dots)$$
$$((\text{array } (n_a \dots n_i \dots) \; \widehat{\mathfrak{v}}_a \dots) \; (\text{Shp } n_i \dots)) \dots$$
$$\iota_r)$$
$$\mapsto_{lift}$$
$$((\text{array } (n_p \dots) \; Concat \left[\!\!\left[ Rep_{n_{fe}} \left[\!\!\left[ Split_1 \left[\!\!\left[ \widehat{\mathfrak{v}}_f \dots \right]\!\!\right] \right]\!\!\right] \right]\!\!\right])$$
$$((\text{array } (n_p \dots \; n_i \dots) \; Concat \left[\!\!\left[ Rep_{n_{ae}} \left[\!\!\left[ Split_{n_{ac}} \left[\!\!\left[ \widehat{\mathfrak{v}}_a \dots \right]\!\!\right] \right]\!\!\right] \right]\!\!\right]) \; (\text{Shp } n_i \dots)) \dots$$
$$\iota_r)$$
where
Not all of $(n_f \dots), (n_a \dots) \dots$ are equal

$$
\begin{aligned}
n_p \dots &= \bigsqcup \left[\!\!\left[ (n_f \dots) \; (n_a \dots) \dots \right]\!\!\right] & n_{fe} &= \prod \frac{(n_p \dots)}{(n_f \dots)} \\
n_{ae} \dots &= \prod \frac{(n_p \dots)}{(n_a \dots)} \dots & n_{ac} \dots &= \left( \prod (n_i \dots) \right) \dots
\end{aligned}
$$

$$((\text{array } (n_f \dots) \; \widehat{\mathfrak{v}}_f \dots)$$
$$((\text{array } (n_f \dots n_i \dots) \; \widehat{\mathfrak{v}}_a \dots) \; (\text{Shp } n_i \dots)) \dots$$
$$\iota_r)$$
$$\mapsto_{map}$$
$$(\text{frame } \iota_r$$
$$((\text{array } () \; \widehat{\mathfrak{v}}_f) \; ((\text{array } (n_i \dots) \; \widehat{\mathfrak{v}}_c \dots) \; (\text{Shp } n_i \dots)) \; \iota_c) \dots)$$
where
$$
\begin{aligned}
n_c \dots &= \left( \prod n_i \dots \right) \dots \\
((\mathfrak{v}_c \dots) \dots) \dots &= Transpose \left[\!\!\left[ Split_{n_c} \left[\!\!\left[ \mathfrak{v}_a \dots \right]\!\!\right] \dots \right]\!\!\right] \\
Length \left[\!\!\left[ n_f \dots \right]\!\!\right] &> 0 \\
\iota_c \dots &= \left( \iota_r \dotdiv (\text{Shp } n_f \dots) \right) \dots
\end{aligned}
$$

$$((\text{array } () \; (\lambda \; (x \dots) \; \widehat{e})) \; ((\text{array } (n_i \dots) \; \widehat{v}) \; (\text{Shp } n_i \dots)) \dots \; \iota_r)$$
$$\mapsto_{\beta} \; \widehat{e}[x \mapsto \widehat{v}, \dots]$$

$$(\text{i-app } (\text{array } (n_f \dots) \; (\text{i}\lambda \; (x \dots) \; \widehat{e}) \; \dots) \; \iota_a \dots \; \iota_r)$$
$$\mapsto_{i\beta} \; (\text{frame } \iota_r \; \widehat{e}[x \mapsto \iota_a, \dots] \dots)$$

$$(\text{frame } (\text{Shp } n \dots) \; (\text{array } (n' \dots) \; \mathfrak{v} \; \dots) \; \dots)$$
$$\mapsto_{collapse} \; (\text{array } (n \dots n' \dots) \; Concat \left[\!\!\left[ (\mathfrak{v} \dots) \dots \right]\!\!\right])$$

$$(\text{unbox } (x_i \dots \; x_e \; (\text{array } (n_s \dots) \; (\text{box } \iota \dots \; \widehat{v}))) \; \widehat{e} \; \iota_b)$$
$$\mapsto_{unbox} \; (\text{frame } (\text{++ } (\text{Shp } n_s \dots) \; \iota_b) \; e[x_i \mapsto \iota, \dots, x_e \mapsto \widehat{v}])$$

Fig. 12. Dynamic semantics for erased Remora

its shape. All occurrences of s from the original program now stand for an array shape rather than a full array type.

We develop a bisimulation argument to show that the behavior of an explicitly typed term matches the behavior of its erased form. We define the space $S$ of machine states to be the sum of the set of well-typed explicit Remora terms and the set of their type-erased forms. That is, $S = Expr^T \uplus \widehat{Expr^T}$, where $Expr^T = \{e \in Expr \,|\, \cdot; \cdot; \cdot \vdash e : \tau\}$ and $\widehat{Expr^T} = \{\mathscr{E}[\![e]\!] \,|\, e \in Expr^T\}$. Transitions in the machine match the explicit and erased languages' respective $\mapsto$ relations. We also define the "erasure equivalence" relation $\cong_{\mathscr{E}}$ on machine

$$\mathscr{E}[\![(\texttt{array}\ (n \ldots)\ \mathfrak{a} \ldots)^{\tau_r}]\!] = (\texttt{array}\ (n \ldots)\ \mathscr{A}[\![\mathfrak{a}]\!] \ldots)$$

$$\mathscr{E}[\![(\texttt{frame}\ (n \ldots)\ e \ldots)^{\tau_r}]\!] = (\texttt{frame}\ (\mathscr{T}[\![\tau_r]\!])\ \mathscr{E}[\![e]\!] \ldots)$$

$$\mathscr{E}\Big[\!\Big[(e_f^{(\texttt{Arr}\ (\texttt{->}\ (\tau_i \ldots)\ \tau_o)\ \iota_f)}\ e_a \ldots)^{\tau_r}\Big]\!\Big] = (\mathscr{E}[\![e_f]\!]\ (\mathscr{E}[\![e_a]\!]\ \mathscr{T}[\![\tau_i]\!]) \ldots \mathscr{T}[\![\tau_r]\!])$$

$$\mathscr{E}[\![(\texttt{t-app}\ e_f\ \tau_a \ldots)^{\tau_r}]\!] = (\texttt{i-app}\ \mathscr{E}[\![e_f]\!]\ \mathscr{T}[\![\tau_a]\!] \ldots \mathscr{T}[\![\tau_r]\!])$$

$$\mathscr{E}[\![(\texttt{i-app}\ e_f\ \iota_a \ldots)^{\tau_r}]\!] = (\texttt{i-app}\ \mathscr{E}[\![e_f]\!]\ \iota_a \ldots \mathscr{T}[\![\tau_r]\!])$$

$$\mathscr{E}[\![(\texttt{unbox}\ (x_i \ldots x_e\ e_s)\ e_b^{\tau_b})]\!] = (\texttt{unbox}\ (x_i \ldots x_e\ \mathscr{E}[\![e_s]\!])\ \mathscr{E}[\![e_b]\!]\ \mathscr{T}[\![\tau_b]\!])$$

$$\mathscr{A}[\![\mathfrak{o}]\!] = \mathfrak{o}$$

$$\mathscr{A}[\![\mathfrak{b}]\!] = \mathfrak{b}$$

$$\mathscr{A}[\![(\lambda\ ((x\ \tau) \ldots)\ e)]\!] = (\lambda\ (x \ldots)\ \mathscr{E}[\![e]\!])$$

$$\mathscr{A}[\![(\texttt{T}\lambda\ ((x\ k) \ldots)\ v)]\!] = (\texttt{I}\lambda\ (x \ldots)\ \mathscr{E}[\![v]\!])$$

$$\mathscr{A}[\![(\texttt{I}\lambda\ ((x\ \gamma) \ldots)\ v)]\!] = (\texttt{I}\lambda\ (x \ldots)\ \mathscr{E}[\![v]\!])$$

$$\mathscr{A}[\![(\texttt{box}\ \iota \ldots e\ \tau)]\!] = (\texttt{box}\ \iota \ldots \mathscr{E}[\![e]\!])$$

$$\mathscr{T}[\![x]\!] = x$$

$$\mathscr{T}[\![(\texttt{Arr}\ \tau\ \iota)]\!] = \iota$$

$$\mathscr{T}[\![\tau]\!] = (\texttt{Shp}) \quad \text{otherwise}$$

$$\mathscr{R}[\![\mathfrak{a}]\!] = \mathscr{A}[\![\mathfrak{a}]\!] \qquad \mathscr{R}[\![e]\!] = \mathscr{E}[\![e]\!]$$

Fig. 13. Type erasure for Remora

states as the equivalence closure of the relation imposed by $\mathscr{E}[\![\cdot]\!]$. Before we show that $\cong_{\mathscr{E}}$ is a bisimulation, several intermediate results are needed.

First, the bisimulation proof will in one case need to reach deep into an expression to find the next redex. A compositionality property of the erasure rule will make it possible to reason about the redex and its reduced form separately from the evaluation context in which it is embedded.

*Lemma 3.1* (*Erasure in context*)
Given an evaluation context $\mathbb{E}$ and expression $e$, where $\mathbb{E}[e]$ is well-typed, $\mathscr{E}[\![\mathbb{E}[e]]\!] = \mathscr{C}[\![\mathbb{E}]\!][\mathscr{E}[\![e]\!]]$.

*Proof*
This follows from straightforward induction on $\mathbb{E}$.    □

We will also rely on a series of lemmas showing that substitution commutes with erasure.

*Lemma 3.2*
$\mathscr{R}[\![t[x \mapsto \mathscr{E}[\![e_x]\!]]]\!] = \mathscr{R}[\![t]\!][x \mapsto \mathscr{E}[\![e_x]\!]]$

*Proof*
This is straightforward induction on $t$.    □

$$\mathscr{C}[\![\Box]\!] = \Box$$

$$\mathscr{C}[\![(\texttt{array}\ (n\dots)\ \mathfrak{v}\dots\ (\texttt{box}\ \iota\dots\ \mathbb{E}\ \tau)\ \mathfrak{a}\dots)]\!]$$
$$= (\texttt{array}\ (n\dots)\ \mathscr{A}[\![\mathfrak{v}]\!]\dots\ (\texttt{box}\ \iota\dots\ \mathscr{C}[\![\mathbb{E}]\!])\ \mathscr{A}[\![\mathfrak{a}]\!]\dots)$$

$$\mathscr{C}[\![(\texttt{frame}\ (n\dots)\ v\dots\ \mathbb{E}\ e\dots)^{\tau_r}]\!]$$
$$= (\texttt{frame}\ (\mathscr{T}[\![\tau_r]\!])\ \mathscr{E}[\![v]\!]\dots\ \mathscr{C}[\![\mathbb{E}]\!]\ \mathscr{E}[\![e]\!]\dots)$$

$$\mathscr{C}\Big[\!\!\Big[(\mathbb{E}^{(\texttt{Arr}\ (\texttt{->}\ (\tau_i\dots)\ \tau_o)\ \iota_f)}\ e_a\dots)^{\tau_r}\Big]\!\!\Big]$$
$$= (\mathscr{C}[\![\mathbb{E}]\!]\ (\mathscr{E}[\![e_a]\!]\ \mathscr{T}[\![\tau_i]\!])\dots\ \mathscr{T}[\![\tau_r]\!])$$

$$\mathscr{C}\Big[\!\!\Big[(e_f^{(\texttt{Arr}\ (\texttt{->}\ (\tau_1\dots\ \tau_2\ \tau_3\dots)\ \tau_o)\ \iota_f)}\ v_1\dots\ \mathbb{E}\ e_3\dots)^{\tau_r}\Big]\!\!\Big]$$
$$= (\mathscr{E}[\![e_f]\!]\ (\mathscr{E}[\![v_a]\!]\ \mathscr{T}[\![\tau_1]\!])\dots\ (\mathscr{C}[\![\mathbb{E}]\!]\ \mathscr{T}[\![\tau_2]\!])\ (\mathscr{E}[\![e_3]\!]\ \mathscr{T}[\![\tau_3]\!])\dots)$$

$$\mathscr{C}[\![(\texttt{t-app}\ \mathbb{E}\ \tau_a\dots)^{\tau_r}]\!]$$
$$= (\texttt{i-app}\ \mathscr{C}[\![\mathbb{E}]\!]\ \mathscr{T}[\![\tau_a]\!]\dots\ \mathscr{T}[\![\tau_r]\!])$$

$$\mathscr{C}[\![(\texttt{i-app}\ \mathbb{E}\ \iota_a\dots)^{\tau_r}]\!]$$
$$= (\texttt{i-app}\ \mathscr{C}[\![\mathbb{E}]\!]\ \iota_a\dots\ \mathscr{T}[\![\tau_r]\!])$$

$$\mathscr{C}[\![(\texttt{unbox}\ (x_i\dots\ x_e\ \mathbb{E})\ e_b)]\!]$$
$$= (\texttt{unbox}\ (x_i\dots\ x_e\ \mathscr{C}[\![\mathbb{E}]\!])\ \mathscr{E}[\![e_b]\!])$$

Fig. 14. Type-erasing Remora evaluation contexts

*Lemma 3.3*
$$\mathscr{T}[\![\tau[x \mapsto \tau_x]]\!] = \mathscr{T}[\![\tau]\!][x \mapsto \mathscr{T}[\![\tau_x]\!]]$$

*Proof*
This is straightforward induction on $\tau$.    □

*Lemma 3.4*
$$\mathscr{R}[\![t[x \mapsto \tau_x]]\!] = \mathscr{R}[\![t]\!][x \mapsto \mathscr{T}[\![\tau_x]\!]]$$

*Proof*
This is straightforward induction on $t$.    □

*Lemma 3.5*
$$\mathscr{T}[\![\tau[x \mapsto \iota_x]]\!] = \mathscr{T}[\![\tau]\!][x \mapsto \iota_x]$$

*Proof*
This is straightforward induction on $\tau$.    □

*Lemma 3.6*
$$\mathscr{R}[\![t[x \mapsto \iota_x]]\!] = \mathscr{R}[\![t]\!][x \mapsto \iota_x]$$

*Proof*
This is straightforward induction on $t$.    □

*Lemma 3.7* (*Values erase to values*)
For any well-typed term $t$,

- If $t$ has the form $\mathfrak{v}$, then $\mathscr{R}[\![t]\!]$ has the form $\widehat{\mathfrak{v}}$

- If $t$ has the form $v$, then $\mathscr{R}[\![t]\!]$ has the form $\widehat{v}$

*Proof*

We use induction on $t$. The only nontrivial cases are `box` and `array` forms, which may be values or may contain incomplete computation. The contents of a `box` value must itself be a value, which the induction hypothesis implies will erase to a value. Similarly, an `array` value contains only atomic values, which erase to atomic values. □

*Lemma 3.8* (*Lockstep*)

For any well-typed $e$, one of the following holds:

- $e$ has the form $v$, and $\mathscr{E}[\![e]\!]$ has the form $\widehat{v}$
- $e \mapsto e'$, and $\mathscr{E}[\![e]\!] \mapsto \mathscr{E}[\![e']\!]$
- $e \not\mapsto$, and $\mathscr{E}[\![e]\!] \not\mapsto$

*Proof*

We prove this by induction on $e$.

We rely on Lemma 3.1 (erasure in context) when $e$ is a redex $e_r$ within an evaluation context $\mathbb{E}$ other than $\square$. If $e_r \not\mapsto$ because we have a mis-applied primitive operator, then the same is true for $\mathscr{E}[\![e_r]\!]$, so $\mathscr{E}[\![e]\!]$ is also an evaluation context around a mis-applied primitive operator. Otherwise, $e_r \mapsto e'_r$, and the induction hypothesis implies that $\mathscr{E}[\![e_r]\!] \mapsto \mathscr{E}[\![e'_r]\!]$. So $\mathscr{E}[\![e]\!] \mapsto \mathscr{E}[\![e']\!]$, the erased context filled with $e'_r$.

Values are handled by Lemma 3.7. For the remaining cases—redexes—straightforward symbol pushing shows that erased Remora's reduction rules follow those of Remora. □

Since we have a deterministic operational semantics for both explicitly typed Remora and type-erased Remora, the lockstep lemma also works in reverse. If an erased term takes an evaluation step, its preimage cannot be a value form or stuck state. The preimage must therefore step to some result expression, which itself erases to the same result. Similarly, a value form or stuck state in erased Remora cannot have a preimage which takes an evaluation step.

*Corollary 3.1* (*Reverse lockstep*)

If $\mathscr{E}[\![e]\!] \mapsto \mathscr{E}[\![e']\!]$, then for any $e''$ such that $e \mapsto e''$, we have $e' \cong_{\mathscr{E}} e''$, and at least one such $e''$ exists. If $\mathscr{E}[\![e]\!] \not\mapsto$, then $e \not\mapsto$.

Recall our relation $\cong_{\mathscr{E}}$ on the set of machine states $S = Expr^T \uplus \widehat{Expr^T}$, where $Expr^T = \{e \in Expr | \cdot; \cdot; \cdot \vdash e : \tau\}$, *i.e.* the set of well-typed explicitly-typed terms, and $\widehat{Expr^T}$ is the image of $Expr^T$ under type erasure. $\cong_{\mathscr{E}}$ is the equivalence closure of the relation given by the erasure function $\mathscr{E}[\![\cdot]\!]$. That is, $\cong_{\mathscr{E}}$ is the least relation which relates two states $s$ and $w$ iff any of the following hold:

1. $s \in Expr^T$ and $\mathscr{E}[\![s]\!] = w$ (erasure proper)
2. $s =_\alpha w$ (reflexivity)
3. $w \cong_{\mathscr{E}} s$ (symmetry)
4. $s \cong_{\mathscr{E}} s'$ and $s' \cong_{\mathscr{E}} w$ (transitivity)

Expanding the erasure relation based on $\mathscr{E}[\![\cdot]\!]$ to include both symmetry and transitivity relates any two explicitly typed expressions which produce $\alpha$-equivalent erased terms. A $\cong_{\mathscr{E}}$

equivalence class consists of a single erased Remora expression and all of its preimages. There can be only one erased Remora expression because type erasure is a well-defined function (*i.e.*, no single explicitly typed expression can erase to multiple different results). Formally, every $\cong_{\mathscr{E}}$ equivalence class must have the form

$$\{\widehat{e}\} \uplus \{e \in Expr \,|\, \mathscr{E}[\![e]\!] = \widehat{e}\}$$

*Theorem 3.1*
$\cong_{\mathscr{E}}$ is a bisimulation. That is, for any states $s, w \in S$ if $s \cong_{\mathscr{E}} w$, either $(s \mapsto u \wedge w \mapsto v \wedge u \cong_{\mathscr{E}} v)$ or $(s \not\mapsto \wedge w \not\mapsto)$.

*Proof*
There are four cases to consider, depending on which of *Expr* or $\widehat{Expr}$ each related term is drawn from, but we can merge the two cases where $s$ and $w$ are drawn from different languages.

$s \in \widehat{Expr}$ **and** $w \in Expr$**, or vice versa:** Then $s$ is the sole type-erased expression in its equivalence class, and $\mathscr{E}[\![w]\!] = s$ (or vice versa). Our proof obligation is exactly the lockstep lemma (Lemma 3.8).

$s, w \in \widehat{Expr}$**:** Since each equivalence class contains only one type-erased expression, $s = w$. They must therefore have the same reduction behavior.

$s, w \in Expr$**:** If $s \not\mapsto$, then the lockstep lemma implies $\mathscr{E}[\![s]\!] = \mathscr{E}[\![w]\!] \not\mapsto$. Then by reverse lockstep, $w \not\mapsto$ as well. On the other hand, if $s \mapsto s'$, then $\mathscr{E}[\![s]\!] = \mathscr{E}[\![w]\!] \mapsto \mathscr{E}[\![s']\!]$. Lockstep implies $\mathscr{E}[\![w]\!] \mapsto \mathscr{E}[\![w']\!]$. Since Erased Remora has deterministic operational semantics, $\mathscr{E}[\![s']\!] = \mathscr{E}[\![w']\!]$ (they are both the result of taking an evaluation step from the same expression). Therefore, all of their preimages, including $s'$ and $w'$ are erasure-equivalent, *i.e.*, $s' \cong_{\mathscr{E}} w'$. $\quad\square$

# 4 Related Work

Rank polymorphism originally appeared in APL (Iverson, 1962), which Iverson designed as a form of mathematical notation, with the APL interpreter serving to eliminate the semantic ambiguity found in conventional notation (Iverson, 1980). At first, APL only lifted *scalar* functions to operate on aggregate data via pointwise application, either on a scalar argument and an aggregate argument or on two aggregates of identical shape. Subsequent development introduced the notion of function rank, the number of dimensions a function expects its argument to have. This generalized the scalar function lifting, *e.g.*, allowing a vector-mean function to produce a vector of results when given a matrix argument, and it introduced the "frame of cells" view of aggregate arguments where pointwise lifting generalizes to cellwise lifting. The next generalization step was to loosen the rule on frame shape compatibility. In J (Jsoftware, Inc., n.d.), which Iverson created as a successor to APL, a function can be applied to two arguments of differing frames as long as one frame— viewed as a sequence of dimensions—is a prefix of the other. This was a conscious design

decision on Iverson's part: prefix agreement was chosen over suffix agreement because it fit better with APL's emphasis on operating along arrays' major axes (Hui, 1995).

FISh also made implicit aggregate lifting part of the semantics of function application, and its static semantics resolves the shapes of all arrays (Jay, 1998). A conventional type judgment describes the elements of arrays computed by a FISh program, and a second judgment ascribes a shape to each array. In FISh's metatheory, the shape of a function is a function on shapes. Thus a function application's shape can be calculated statically by applying the shape function to the arguments' shapes.

As elegant as this model is in a first-order language, it is incompatible with first-class functions. When functions appear in arrays which are themselves applied to arguments, the shape must describe the layout of of that collection of functions, not just the functions' own behavior. In Remora, a function which checks whether a point in $\mathbb{R}^3$ is near the origin might have the type `(-> ((Arr Float (Shp 3))) (Arr Bool (Shp)))`. FISh considers this function's shape to be the function on a singleton domain (containing only [3]) which returns the empty vector. However, Remora expressions produce array data, even in function position. A function array `near-origin?` containing only that function has its own "first-order" shape `(Shp)`, independent of any function on shapes summarizing its behavior.

In resolving all shapes statically, FISh is also too restrictive to permit shapes determined from run-time data. Functions like `iota` and `filter` cannot exist, nor can the ragged data which would result from lifting them. By characterizing functions with restricted dependent types, Remora escapes both of these limitations.

ZPL is a data-parallel language which was designed to live within a larger language with a more general parallelism mode (Lin & Snyder, 1994). Its programming model is based on an explicit map operation over programmer-specified index space within an array. Several built-in operators modify an index space, such as shifting a section along some dimension, adding a new dimension by broadcasting, or slicing out a particular sub-array. The set of built in operators is constructed to make communication cost implications clear to the programmer (Chamberlain *et al.*, 1998).

NESL uses a nested-vector data model, rather than the rank-polymorphic regular-array model (Blelloch, 1995). Programs map operations over vectors using a comprehension notation. Since a vector's elements may themselves be vectors with widely varying lengths, NESL's main performance trick is turning irregular nested vectors into a flat internal representation. After flattening, NESL can evenly distribute the computation workload by splitting at places that user code doe not consider not sub-vector boundaries.

The goal of isolating the cell-level portion of a program has also led to some domain-specific languages, targeting specific varieties of regular aggregate data. Halide is a language designed for writing image processing pipelines (Ragan-Kelley *et al.*, 2012). The iteration space is the set of pixels in an image. A Halide programmer writes code describing how an individual pixel should be handled (or a cluster of pixels for stencil computation), and then in a separate portion of the program, the programmer describes how the iteration space ought to be traversed. Since the program does not intermingle loop-nest control code with loop-body computation code, it is easier for a maintainer to tune for performance by adjusting the iteration schedule.

Diderot is a programming language specifically aimed at processing medical images (Chiw *et al.*, 2012). The universe of data consists of tensor fields, intended as functions on the continuous domain $\mathbb{R}^n$, rather than any particular discrete index space. Diderot offers pointwise arithmetic on tensors and a collection of common operations such as outer product and transposition. Aggregate lifting appears again, albeit in a smaller form, with some operations on tensor fields.

HorseIR uses implicit lifting over arrays as an intermediate representation for SQL queries (Chen *et al.*, 2018). HorseIR is lower-level than APL itself, serving as a three-address IR with vector instructions, and the IR itself is designed to ease loop fusion. In contrast with APL, all HorseIR arrays are vectors—there are no higher-rank arrays, and scalars are represented as unit-length vectors. There is also a list datatype for handling heterogeneous aggregate data, such as one row of a database table.

Our type system's use of restricted dependent types is inspired by Dependent ML (Xi, 1998). While Dependent ML is designed with the expectation that the index language has a fully decidable theory, Remora's index language does not (Durnev, 1995). In subsequent work, Dependent ML also focused on ensuring safety of array index accesses, using singleton and range types to ensure that numbers used for indexing fell within arrays' bounds (Xi & Pfenning, 1998).

Others have applied established type system machinery to an APL-like computation model. Thatte's coercion semantics (Thatte, 1991) uses a form subtyping in which scalar types are subtypes of aggregates, and aggregate types in certain situations are subtypes of higher-dimensional aggregates. The subtyping rules emit coercions which invoke functions such as `map` and `replicate`, automatically adapting scalar functions to aggregate data. However, the restrictions on treating aggregate types as subtypes of larger aggregates prohibit lifting for functions which expect non-scalar input (*i.e.*, those whose expected rank is greater than zero) and lifting to unequal frame shapes (*e.g.*, vector-matrix addition). Gibbons's embedding of in an extended Haskell (Gibbons, 2016) extensively uses type-level programming and defines much of the rank-polymorphic lifting machinery in terms of transposition.

Single Assignment C (SAC) is a variant of C without mutation (Scholz, 2003). The primary iteration mechanism in SAC is the `with` loop, in which the programmer describes a traversal of the index space and builds an output array an element at a time. The lack of mutation pushes the programmer to avoid writing loop-carried dependence, much like in the rank-polymorphic model. Translation of well-behaved APL programs to SAC tends to be straightforward (Grelck & Scholz, 1998), though the resulting loop structure is somewhat specific to the function being lifted. A SAC variant called Qube introduces a Dependent ML-style type system to ensure the safety of indexed array element accesses (Trojahner & Grelck, 2009). Qube retains SAC's `with` loop-based programming model—relying on range types in contrast to APL's and Remora's implicitly lifted function application semantics—and due to its C roots, it does not support first-class functions to the extent of permitting arrays to appear in function position.

Since the syntactic structure of a line of APL code, in particular the meaning of the juxtaposition of two terms, depends on the meanings of names which appear in the code, standard APL does not admit a fully static parsing algorithm. Due to this and other idiosyn-

cratic warts, past efforts to compile APL have typically targetted large but well-behaved subsets of the language.

A prominent exception is the APL\3000 compiler, which could produce machine code for individual statements and used interpretation to manage inter-statement control flow (Johnston, 1979). This allowed some intraprocedural optimization, such as fold-unfold and map-map fusion[3] (Abrams, 1970), analogous to deforestation which arose in mainstream functional languages (Burstall & Darlington, 1977; Wadler, 1984). It also ensured that names' dynamic meanings would be available by the time execution reached any statement which used them.

Weiss and Saal instead applied interprocedural data-flow analysis to resolve the syntactic classes of variable names in APL code (Weiss & Saal, 1981). This analysis is not complete for APL itself due to the possibility that reassignment of a variable name will change how some line of code parses. However, the authors found that real-world code did not make use of the full freedom to manipulate the syntactic structure of an APL program by dynamically reassigning variables, suggesting that this is a misfeature which can be discarded for little cost. However, the common line-at-a-time compilation style used in other work also served to ensure that by the time execution reaches a line of code, its variable names, and thus their syntactic classes, could be resolved.

Budd describes a compiler for an APL variant in which the ambiguity in parsing is avoided by declaring identifiers before use and name resolution is simplified by adopting lexical scope. (Budd, 1988).

APEX parallelizes an APL dialect with restrictions on many APL features deemed incompatible with compilation (such as producing a string representation of an arbitrary function) or not strictly necessary for practical use (such as `GOTO`) (Bernecky, 1999). APEX's shape compatibility rules for implicit aggregate lifting are less permissive than APL, *e.g.*, prohibiting most cases of vector-matrix addition. Instead, both arguments must have the same rank or at least one argument passed to a scalar-consuming function be a scalar or singleton vector.

A more recent line of work has focused on intermediate representations of array programs, such as the Typed Array Intermediate Language, which makes aggregate operations explicit using an `each` primitive operator (Elsman & Dybdal, 2014). The type system tracks arrays' ranks, which provides enough information to recognize when an `each` call is needed, but it is not meant to ensure that lifting a function application has a well-defined result (*i.e.*, shape incompatibility is still possible). $\mathscr{L}_0$ plots out a loop fusion strategy inspired by control flow graph reduction (Henriksen & Oancea, 2013). The $\mathscr{L}_0$ compiler splits array programs into kernels of fused aggregate operations based on the applicability of $T_2$ graph reductions (merging nodes $X$ and $Y$ if $X$ is $Y$'s sole predecessor) to the program's data flow graph.

---

[3] In the APL community, these transformations are respectively referred to as "beating" (imagine a metronome producing a stream of beats) and "dragging along" (amassing a chain of delayed scalar operations to apply in a single loop body).

## 5 Conclusion

While APL and its descendant languages have attracted a devoted userbase, there has been little cross talk between the array-language community and the broader programming language research community. Much of the analysis opportunity taken for granted by lambda calculists has been unavailable in APL—despite the many implementations, such languages lacked formal semantics amenable to proofs. Meanwhile, implementations of rank-polymorphic languages have struggled with compilation due to control structure that is "too dynamic," depending on run-time data to determine the structure of a loop nest.

Developing Remora's semantics kills two birds with one stone: Formally stated reduction rules describe the results expected from the implicit, data-driven control structure, and typing rules give enough information about array shapes to identify that control structure statically. This necessarily entails recognizing programs which cannot have such a control structure due to incompatible array shapes—this is not what we set out to do but a benefit realized by pursuing a larger goal. By casting the aggregate lifting as an extension to $\lambda$-calculus's function application semantics, we escape from APL's limitation on function arity and can treat functions as first-class values.

There are two things to look for in evaluating a type system. We want to know the conclusions drawn by type checking reflect reality, which is handled by a conventional type soundness theorem. We also want to be sure that the types convey the information we seek. In Remora's case, that information is the iteration structure implicit in each function application. The typing rule for function application (and similarly, the rules for type and index application) produces a static characterization of that implicit iteration structure.

Our type erasure can be seen as a compilation pass which moves the decision about how to break arguments into cells from the function's type into the application term. While our primary purpose in presenting type erasure is to point out type-level information which is not truly needed at run time, it also serves to make the control structure one step more explicit. Arguments' full shapes—available both at erasure time and later by inspecting argument terms—suffice to determine the arguments' frames, the last puzzle piece needed to turn Remora's implicit iteration structure into explicit calls to `map` and `replicate` functions. Our intention for future work is to demonstrate use of that shape information for fully static compilation of Remora code.

## References

Abrams, Philip S. (1975). What's wrong with APL? *Pages 1–8 of: Proceedings of Seventh International Conference on APL.* APL '75. New York, NY, USA: ACM.

Abrams, Philip Samuel. (1970). *An apl machine.* Ph.D. thesis, Stanford, CA, USA. AAI7022146.

Backus, John. (1978). Can programming be liberated from the Von Neumann style?: A functional style and its algebra of programs. *Commun. ACM*, **21**(8), 613–641.

Bernecky, Robert. (1999). *APEX, the APL parallel executor.* M.Phil. thesis, National Library of Canada= Bibliothèque nationale du Canada.

Blelloch, Guy. (1995). *NESL: A nested data-parallel language (version 3.1).* Tech. rept.

Budd, Timothy. (1988). *An APL compiler.* Springer-Verlag.

Burstall, R. M., & Darlington, John. (1977). A transformation system for developing recursive programs. *J. acm*, **24**(1), 44–67.

Chamberlain, B., Choi, S., Lewis, E., Lin, C., Snyder, L., & Weathersby, W. (1998). ZPL's WYSIWYG performance model. *Pages 50– of: Proceedings of the High-Level Parallel Programming Models and Supportive Environments*. HIPS '98. Washington, DC, USA: IEEE Computer Society.

Chen, Hanfeng, D&#39;silva, Joseph Vinish, Chen, Hongji, Kemme, Bettina, & Hendren, Laurie. (2018). Horseir: Bringing array programming languages together with database query processing. *Pages 37–49 of: Proceedings of the 14th ACM SIGPLAN International Symposium on Dynamic Languages*. DLS 2018. New York, NY, USA: ACM.

Chiw, Charisee, Kindlmann, Gordon, Reppy, John, Samuels, Lamont, & Seltzer, Nick. (2012). Diderot: A parallel DSL for image analysis and visualization. *Pages 111–120 of: Proceedings of the 33rd ACM SIGPLAN Conference on Programming Language Design and Implementation*. PLDI '12. New York, NY, USA: ACM.

Durnev, Valery G. (1995). Undecidability of the positive $\forall\exists^3$-theory of a free semigroup. *Siberian mathematical journal*, **36**(5), 917–929.

Elsman, Martin, & Dybdal, Martin. (2014). Compiling a subset of APL into a typed intermediate language. *Pages 101:101–101:106 of: Proceedings of ACM SIGPLAN International Workshop on Libraries, Languages, and Compilers for Array Programming*. ARRAY'14. New York, NY, USA: ACM.

Gibbons, Jeremy. (2016). APLicative programming with naperian functors (extended abstract). *Pages 13–14 of: Proceedings of the 1st International Workshop on Type-Driven Development*. TyDe 2016. New York, NY, USA: ACM.

Grelck, Clemens, & Scholz, Sven-Bodo. (1998). Accelerating APL programs with SAC. vol. 29. New York, NY, USA: ACM.

Henriksen, Troels, & Oancea, Cosmin Eugen. (2013). A t2 graph-reduction approach to fusion. *Pages 47–58 of: Proceedings of the 2Nd ACM SIGPLAN Workshop on Functional High-performance Computing*. FHPC '13. New York, NY, USA: ACM.

Hui, Roger K. W. (1995). Rank and uniformity. *SIGAPL APL quote quad*, **25**(4), 83–90.

Iverson, Kenneth E. (1962). *A programming language*. New York, NY, USA: John Wiley & Sons, Inc.

Iverson, Kenneth E. (1980). Notation as a tool of thought. *Commun. ACM*, **23**(8), 444–465.

Jay, C. B. (1998). *The FISh language definition*. Tech. rept.

Johnston, Ronald L. (1979). The dynamic incremental compiler of APL\3000. *Pages 82–87 of: Proceedings of the International Conference on APL: Part 1*. APL '79. New York, NY, USA: ACM.

Jsoftware, Inc. *Jsoftware: High-performance development platform*.

Lin, Calvin, & Snyder, Lawrence. (1994). ZPL: An array sublanguage. *Pages 96–114 of: Proceedings of the 6th International Workshop on Languages and Compilers for Parallel Computing*. London, UK, UK: Springer-Verlag.

Mathworks, Inc. (1992). *Matlab, user's guide*. Natick, MA.

Oliphant, Travis E. 2006 (12). *Guide to numpy*. Trelgol Publishing USA.

Ragan-Kelley, Jonathan, Adams, Andrew, Paris, Sylvain, Levoy, Marc, Amarasinghe, Saman, & Durand, Frédo. (2012). Decoupling algorithms from schedules for easy optimization of image processing pipelines. *Acm trans. graph.*, **31**(4), 32:1–32:12.

Scholz, Sven-Bodo. (2003). Single assignment C: efficient support for high-level array operations in a functional setting. *J. funct. program.*, **13**(6), 1005–1059.

Thatte, Satish. (1991). A type system for implicit scaling. *Sci. comput. program.*, **17**(1-3), 217–245.

Trojahner, Kai, & Grelck, Clemens. (2009). Dependently typed array programs don't go wrong. *Journal of logic and algebraic programming*, **78**(7), 643–664.

Wadler, Philip. (1984). Listlessness is better than laziness: Lazy evaluation and garbage collection at compile-time. *Pages 45–52 of: Proceedings of the 1984 ACM Symposium on LISP and Functional Programming*. LFP '84. New York, NY, USA: ACM.

Weiss, Zvi, & Saal, Harry J. (1981). Compile time syntax analysis of APL programs. *Pages 313–320 of: Proceedings of the International Conference on APL*. APL '81. New York, NY, USA: ACM.

Xi, Hongwei. (1998). *Dependent types in practical programming*. Ph.D. thesis, Pittsburgh, PA, USA. AAI9918624.

Xi, Hongwei, & Pfenning, Frank. (1998). Eliminating array bound checking through dependent types. *Pages 249–257 of: Proceedings of the ACM SIGPLAN 1998 Conference on Programming Language Design and Implementation*. PLDI '98. New York, NY, USA: ACM.