

A Simple and Practical Linear Algebra Library Interface with Static Size Checking*

Akinori Abe Eijiro Sumii

Graduate School of Information Sciences
Tohoku University, Japan

abe@sf.ecei.tohoku.ac.jp sumii@sf.ecei.tohoku.ac.jp

Linear algebra is a major field of numerical computation and is widely applied. Most linear algebra libraries (in most programming languages) do not statically guarantee consistency of the dimensions of vectors and matrices, causing runtime errors. While advanced type systems—specifically, dependent types on natural numbers—can ensure consistency among the sizes of collections such as lists and arrays [3, 4, 22], such type systems generally require non-trivial changes to existing languages and application programs, or tricky type-level programming.

We have developed a linear algebra library interface that verifies the consistency (with respect to dimensions) of matrix operations by means of *generative phantom types*, implemented via fairly standard ML types and module system. To evaluate its usability, we ported to it a practical machine learning library from a traditional linear algebra library. We found that most of the changes required for the porting could be made mechanically, and changes that needed human thought are minor.

1 Introduction

Linear algebra is a major field of numerical computation and is widely applied to many domains such as image processing, signal processing, and machine learning. Some programming languages support built-in matrix operations, while others are equipped with libraries. Examples of the former are Matlab, statistical programming language S, and OptiML [18, 20] (a domain-specific language for machine learning embedded in Scala). The latter include BLAS [16] and LAPACK [17], originally developed for Fortran and now ported to many languages.

Most of the programming languages and linear algebra libraries do not statically verify consistency of the dimensions of vectors and matrices. Dimensional inconsistency like addition of a 3-by-5 matrix and a 5-by-3 matrix causes runtime errors such as exceptions or, worse, memory corruption.

Advanced type systems—specifically, dependent types on natural numbers—can statically ensure consistency among the sizes of collections such as lists and arrays [3, 4, 22]. However, such type systems generally require non-trivial changes to existing languages and application programs, or tricky type-level programming.

We have developed a linear algebra library interface that guarantees the consistency (with respect to dimensions) of matrix (and vector) operations by using *generative phantom types* as fresh identifiers for statically checking the equality of sizes (i.e., dimensions). This interface has three attractive features in particular.

- It can be implemented only using fairly standard ML types and module system. Indeed, we implemented the interface in OCaml (without significant extensions like GADTs) as a wrapper for an existing library.

*This work was partially supported by JSPS KAKENHI Grant Numbers 22300005, 25540001, 15H02681, and by Mitsubishi Foundation Research Grants in the Natural Sciences.

- For most high-level operations on matrices (e.g., addition and multiplication), the consistency of sizes is verified statically. (Certain low-level operations, like accesses to elements by indices, need dynamic checks.)
- Application programs in a traditional linear algebra library can be easily migrated to our interface. Most of the required changes can be made mechanically.

We implemented our static size checking scheme as a linear algebra library interface—that we call SLAP (Sized Linear Algebra Package, <https://github.com/akabe/slap/>)—on top of an existing linear algebra library LACAML [15]. To evaluate the usability of our interface, we ported to it a practical machine learning library OCaml-GPR [14] from LACAML, thereby ensuring the consistency of sizes.

This paper is structured as follows. In the next section, we explain our basic idea of static size checking through examples. In Section 3, we describe the implementation of our library interface. In Section 4, we report the changes required for the porting of OCaml-GPR along with the number of lines changed. We compare our approach with related work in Section 5 and conclude in Section 6.

2 Our idea

Let `'n vec` be the type of `'n`-dimensional vectors, `('m, 'n) mat` be the type of `'m`-by-`'n` matrices, and `'n size` be the *singleton type* on natural numbers as sizes of vectors and matrices, i.e., evaluation of a term of type `'n size` always results in the natural number corresponding to `'n`. The formal type parameters `'m` and `'n` are instantiated with actual types that represent the sizes of the vectors or matrices. Here we only explain how the dimensions of vectors are represented since those of matrices (as well as sizes) can be represented similarly.

The abstract type `'n vec` can be implemented as any data type that can represent vectors, e.g., `float array`, where the type parameter `'n` is phantom, meaning that it does not appear on the right hand side of the type definition. A phantom type parameter is often instantiated with a type that has no value (i.e., no constructor), which we call a *phantom type*¹. The type `'n vec` must be made abstract by hiding its implementation in the module signature so that the size information in the (phantom) type parameter `'n` is not ignored by the typechecker.

It is relatively straightforward to represent dimensions (size information) as types by, for example, using type-level natural numbers when this information is decided at compile time. The main problem is how to represent dimensions that are unknown until runtime. It is practically important because such dynamically determined dimensions are common (e.g., the length of a vector loaded from a file). Consider the following code for example:

```
let (x : ?1 vec) = loadvec "file1" in
let (y : ?2 vec) = loadvec "file2" in
add x y (* add : 'n vec → 'n vec → 'n vec *)
```

The function `loadvec` of type `string → ? vec` returns a vector of some dimension, loaded from the given path. The third line should be ill-typed because the dimensions of `x` and `y` are probably different. (Even if `"file1"` and `"file2"` were the same path, the addition should be ill-typed because the file might have changed between the two loads.) Thus, the return type of `loadvec` should be different every

¹In [1], the term “phantom type” has the same meaning as in this paper. However, it is used differently in some other papers: types *with* phantom type parameters [13]—or GADT (generalized algebraic data type) [10]—are called “phantom types”.

time it is called (regardless of the specific values of the argument). We call such a return type *generative* because the function returns a value of a fresh type for each call.

The vector type with generative size information essentially corresponds to an existentially quantified sized type like $\exists n. n \text{ vec}$. Our basic idea is to verify *only* the equality of dimensions by representing them as (only) generative phantom types. We implemented this idea in OCaml and carried out a realistic case study to demonstrate that it is mostly straightforward to write (or port) application programs by using our interface (see Section 4 for details).

2.1 Implementation of generative phantom types

We implement the idea of generative phantom types via packages of types like $\exists n. n \text{ vec}$ by using first-class modules and (generative) functors. For convenience, assume that the abstract type $'n \text{ vec}$ is implemented as `float array`.

With first-class modules, the function `loadvec` can be defined as

```

module type VEC = sig
  type n          (* corresponds to a generative phantom type '?' *)
  val value : n vec (* corresponds to '? vec' *)
end

let loadvec filename = (module struct
  type n
  let value = loadarray filename
  end : VEC)

```

where `loadarray : string → float array` loads an array from the given path. The existential type `n` in the module returned by `loadvec` is different every time it is called:

```

module X = (val loadvec "file1" : VEC)
module Y = (val loadvec "file2" : VEC)
add X.value Y.value (* ill-typed since X.n ≠ Y.n *)

```

Alternatively, this behavior can also be implemented by using generative functors² as follows.

```

module Loadvec (X : sig val filename : string end) () : VEC = struct
  type n
  let value = loadarray filename
end

module X = Loadvec(struct let filename = "file1" end) ()
module Y = Loadvec(struct let filename = "file2" end) ()
add X.value Y.value (* ill-typed since X.n ≠ Y.n *)

```

Let us compare the two approaches at the caller site (i.e., from the viewpoint of a user of our interface) in the code above. With first-class modules, the user needed to write the signature `VEC` for unpacking, which was not required for generative functors.

In some cases, however, both approaches require more annotations. Specifically, consider `columns : 'm vec list → ('m, ?) mat`, which creates a matrix by concatenating column vectors in a given list. The number of rows in the matrix returned by `columns` is equal to the dimension of the given

²Generative functors have been supported since OCaml 4.02. Before 4.02, they can be simulated by always passing a module expression of the form `struct . . . end` (rather than a named module like `M`) to applicative functors.

<pre> (* definition: *) let columns (type k) (l : k vec list) = (module struct type n and m = k let value = ... end : MAT with type m = k) (* caller site: *) let f (type k) ... = let l0 = ... in let module X = (val concat l0 : MAT with type m = k) in ... </pre>	<pre> (* definition: *) module Columns (L : sig type k val l : k vec list end) () : MAT with type m = L.k = struct type n and m = L.k let value = ... end (* caller site: *) let f (type k') ... = let l0 = ... in let module X = Concat(struct type k = k' let l = l0 end) () in ... </pre>
---	--

(a) using first-class modules

(b) using a functor

Figure 1: Implementations and callers of `columns`

vectors, while the number of columns is the same as the length of the list, so the latter is generative while the former is not. Figure 1 shows implementations and callers of `columns` using first-class modules and a functor. The signature `MAT` in Figure 1 is defined as:

```

module type MAT = sig
  type m and n
  val value : (m, n) mat
end

```

On one hand, when using first-class modules, the signature `MAT` of the returned module `columns l0` needs the sharing constraint **with type** `m = k`. On the other hand, with functors, the actual type argument **type** `k = k'` cannot be omitted.

In summary, the two approaches require different styles of annotations. We therefore adopted both and often provided two interfaces for the same function so that a user can choose.

2.2 Free type constructors to represent operations on natural numbers

Results of some deterministic functions could also be given generative phantom types but can be given a more precise type using free (in the sense of free algebra) constructors. For example, consider the `append` function, which concatenates an m -dimensional vector and an n -dimensional vector, and returns a vector of dimension $m + n$. Theoretically, its typing would be

```

val append : 'm vec → 'n vec → ('m + 'n) vec

```

if there were addition `+` of type-level natural numbers. Using generative phantom types, however, it can instead be typed as follows

```

val append : 'm vec → 'n vec → ? vec

```

or, using a type constructor `add`,

```
val append : 'm vec → 'n vec → ('m, 'n) add vec
```

The last alternative can be implemented directly in OCaml (or any language with ML-like parametrized data types) with a (preferably phantom) type `('m, 'n) add`, which represents a size that differs from any other (e.g., `('m, 'n) add` differs from `('n, 'm) add` since the constructor `add` is “free”). The return type of `append` does not need to be generative because the dimension of the returned vector is uniquely determined by the dimensions (i.e., the types) of the argument vectors. The same technique can also be applied for other functions such as `cons` (which adds an element to the head of a given vector) and `tl` (which takes a subvector without the head element.)

3 Typing of BLAS and LAPACK functions

BLAS (Basic Linear Algebra Subprograms) [16] and LAPACK (Linear Algebra PACKage) [17] are major linear algebra libraries for Fortran. To evaluate the effectiveness of our idea, we implemented a linear algebra library interface as a “more statically typed” wrapper of LACAML, which is a BLAS and LAPACK binding in OCaml. Our interface is largely similar to LACAML so that existing application programs can be easily ported. Here we explain several techniques required for typing the BLAS and LAPACK functions.

3.1 Function types that depend on flags

3.1.1 Transpose flags for matrices

The BLAS function `gemm` multiplies two general matrices:

```
val gemm : ?beta:num_type → ?c:mat (* C *) →
  ?transa:[ 'N | 'T | 'C ] → ?alpha:num_type → mat (* A *) →
  ?transb:[ 'N | 'T | 'C ] → mat (* B *) → mat (* C *)
```

Basically, it executes $\mathbf{C} := \alpha\mathbf{AB} + \beta\mathbf{C}$. The parameters `transa` and `transb` specify no transpose (`'N`), transpose (`'T`), or conjugate transpose (`'C`) of the matrices \mathbf{A} and \mathbf{B} respectively. For example, if `transa='N` and `transb='T`, then `gemm` executes $\mathbf{C} := \alpha\mathbf{AB}^T + \beta\mathbf{C}$. Thus, the *types* (dimensions) of the matrices change depending on the *values* of the flags (the transpose of an m -by- n matrix is an n -by- m matrix). To implement this behavior, we give each transpose flag a function type that represents the change in types induced by that particular transposition, like:

```
type 'a trans (* = [ 'N | 'T | 'C ] *)
val normal : (('m, 'n) mat → ('m, 'n) mat) trans (* = 'N *)
val trans : (('m, 'n) mat → ('n, 'm) mat) trans (* = 'T *)
val conjtr : (('m, 'n) mat → ('n, 'm) mat) trans (* = 'C *)
```

```
val gemm : ?beta:num_type → ?c:('m, 'n) mat (* C *) →
  transa:(('am, 'ak) mat → ('m, 'k) mat) trans →
  ?alpha:num_type → ('am, 'ak) mat (* A *) →
  transb:(('bk, 'bn) mat → ('k, 'n) mat) trans →
  ('bk, 'bn) mat (* B *) → ('m, 'n) mat (* C *)
```

The arguments `transa` and `transb` are optional in LACAML, but mandatory in our interface because OCaml restricts the types of parameters to those of the default arguments. This is a shortcoming of the typing of optional arguments in OCaml.

3.1.2 Side flags for square matrix multiplication

The BLAS function `symm` multiplies a symmetric matrix **A** by a general matrix **B**:

```
val symm : ?side:[ `L | `R ] → ?beta:num_type →
  ?c:mat (* C *) → ?alpha:num_type → mat (* A *) →
  mat (* B *) → mat (* C *)
```

The parameter `side` specifies the “direction” of the multiplication: `symm` executes $\mathbf{C} := \alpha\mathbf{AB} + \beta\mathbf{C}$ if `side` is ``L`, and $\mathbf{C} := \alpha\mathbf{BA} + \beta\mathbf{C}$ if it is ``R`. If **B** and **C** are m -by- n matrices, **A** is an m -by- m matrix in the former case and n -by- n in the latter case. We implemented these flags as follows:

```
type ('k, 'm, 'n) side (* = [ `L | `R ] *)
val left  : ('m, 'm, 'n) side (* = `L *)
val right : ('n, 'm, 'n) side (* = `R *)
```

The parameter `'k` in type `('k, 'm, 'n) side` corresponds to the dimension of the `'k`-by-`'k` symmetric matrix **A**, and the other parameters `'m` and `'n` correspond to the dimensions of the `'m`-by-`'n` general matrix **B**. When **A** is multiplied from the left to **B** (i.e., like \mathbf{AB}), `'k` is equal to `'m`; therefore, the type of the flag `left` is `('m, 'm, 'n) side`. Conversely, if **A** is right-multiplied to **B** (i.e., like \mathbf{BA}), `'k` is equal to `'n`. Thus, the flag `right` is given the type `('n, 'm, 'n) side`. By using this trick, we can type `symm` as:

```
val symm : side:(('k, 'm, 'n) side → ?beta:num_type →
  ?c:(('m, 'n) mat (* C *) → ?alpha:num_type →
  ('k, 'k) mat (* A *) → ('m, 'n) mat (* B *) → ('m, 'n) mat
```

The same trick can be applied to other square matrix multiplications as well.

3.1.3 Flags that change the size of the results

LAPACK provides routines `gesdd` and `gesvd` for singular value decomposition (SVD), a variant of the eigenvalue problem. It factorizes a given m -by- n matrix **A** as

$$\mathbf{A} = \mathbf{UDV}^\dagger$$

where **U** is an m -by- m unitary matrix, \mathbf{V}^\dagger is the conjugate transpose of an n -by- n unitary matrix **V**, and **D** is an m -by- n matrix with $\min(m, n)$ diagonal elements. The diagonal elements are called singular values (similar to eigenvalues), and columns of **U** and **V** are called left and right singular vectors (similar to eigenvectors), respectively. The singular vectors are sorted according to the corresponding singular values. The challenge is that `gesdd` and `gesvd` store the singular vectors differently to **U**, **V**, or **A** depending on the flags.

We consider `gesdd` first:

```
val gesdd : ?jobz:[ `A | `N | `O | `S ] →
  ?s:vec (* D *) → ?u:mat (* U *) → ?vt:mat (* V† *) → mat (* A *) →
  vec (* D *) * mat (* U *) * mat (* V† *)
```

It computes all singular values but computation of the singular vectors is optional depending on the flags:

- When the flag `jobz` is ``A`, all the left and right singular values are computed and are stored in `u` and `vt`. In this case, the storage `u` and `vt` must be an m -by- m and an n -by- n matrices, respectively.
- If `jobz` is ``S`, only the top $\min(m, n)$ columns in **U** and the top $\min(m, n)$ rows in \mathbf{V}^\dagger are stored in `u` and `vt`, which must be m -by- $\min(m, n)$ and $\min(m, n)$ -by- n , respectively.

- Flag `'O` specifies to overwrite the matrix \mathbf{A} with the top singular vectors as follows:
 - If $m \geq n$, \mathbf{A} is overwritten with the top $\min(m, n)$ columns of \mathbf{U} , and the n rows of \mathbf{V}^\dagger are returned in `vt`; thus `vt` is n -by- n while `u` is not used.
 - If $m < n$, \mathbf{A} is overwritten with the top $\min(m, n)$ rows of \mathbf{V}^\dagger , and the m columns of \mathbf{U} are returned in `u`, so `u` is m -by- m and `vt` is not used.
- For `'N`, no singular vectors are calculated at all (only singular values are returned).

We implement the dependency of the sizes of `u` and `vt` only on the value of `jobz`, ignoring whether $m \geq n$ or $m < n$ in the case of `'O` (i.e., `u` is required to be m -by- m and `vt` to be n -by- n whenever the flag is `'O` even though only one of them is used; in addition, the SLAP-version of `gesdd` returns a value of type `- vec * - mat option * - mat option` instead of `- vec * - mat * - mat` to avoid allocating dummy matrices when \mathbf{U} or \mathbf{V} are not returned).

Based on the ideas above, we defined the SVD job flags and the type of `gesdd` as follows:

```

type ('a, 'b, 'c, 'd, 'e) svd_job (* = [ 'A | 'N | 'O | 'S ] *)

val svd_all      : ('a, 'a,  -,  -,  -) svd_job (* = 'A *)
val svd_top     : ('a,  -, 'a,  -,  -) svd_job (* = 'S *)
val svd_overwrite : ('a,  -,  -, 'a,  -) svd_job (* = 'O *)
val svd_no      : ('a,  -,  -,  -, 'a) svd_job (* = 'N *)

val gesdd : jobz:('u_cols * 'vt_rows,
                  'm * 'n,
                  ('m, 'n) min * ('m, 'n) min,
                  'm * 'n,
                  z * z)
          svd_job →
          ?s: (('m, 'n) min, cnt) vec          (*  $\mathbf{D}$  *) →
          ?u: ('m, 'u_cols, 'u_cd) mat        (*  $\mathbf{U}$  *) →
          ?vt: ('vt_rows, 'n, 'vt_cd) mat     (*  $\mathbf{V}^\dagger$  *) →
          ('m, 'n, 'a_cd) mat                (*  $\mathbf{A}$  *) →
          (('m, 'n) min, 's_cd) vec          (*  $\mathbf{D}$  *) *
          ('m, 'u_cols, 'u_cd) mat option   (*  $\mathbf{U}$  *) *
          ('vt_rows, 'n, 'vt_cd) mat option (*  $\mathbf{V}^\dagger$  *)

```

In the type of `svd_all`, the first type parameter is the same as the second. When it is passed to `jobz`, the OCaml typechecker unifies `'u_cols * 'vt_rows` and `'m * 'n` so that `u` and `vt` have types `('m, 'm) mat` and `('n, 'n) mat`, respectively. Similarly, if `svd_top` is specified, `u` and `vt` are typed as `('m, ('m, 'n) min) mat` and `(('m, 'n) min, 'n) mat`. In the case of `svd_overwrite`, `u` : `('m, 'm) mat` and `vt` : `('n, 'n) mat` are derived. For `svd_no`, `u` and `vt` have types `('m, z) mat` and `(z, 'n) mat`, where `z` is a nullary type constructor representing 0. (These types are virtually equal to `(z, z) mat` because a matrix of `('m, z) mat` or `(z, 'n) mat` has no element at all.)

Second, we consider the type of `gesvd`. The original typing (in LACAML) of `gesvd` is:

```

val gesvd : ?jobu:[ 'A | 'N | 'O | 'S ] → ?jobvt:[ 'A | 'N | 'O | 'S ] →
          ?s:vec (*  $\mathbf{D}$  *) → ?u:mat (*  $\mathbf{U}$  *) → ?vt:mat (*  $\mathbf{V}^\dagger$  *) → mat (*  $\mathbf{A}$  *) →
          vec (*  $\mathbf{D}$  *) * mat (*  $\mathbf{U}$  *) * mat (*  $\mathbf{V}^\dagger$  *)

```

It takes two SVD job flags `jobu` and `jobvt` for the computation of the singular vectors in \mathbf{U} and \mathbf{V}^\dagger . When `jobu` is `'A`, `'S`, or `'N`, all, the top $\min(m, n)$, or no columns in \mathbf{U} are computed respectively. If it

is `'0`, `a` is overwritten with the top $\min(m, n)$ columns. The meaning of `jobvt` is similar. (It is a runtime error to give `'0` for both `jobu` and `jobvt`, because \mathbf{A} cannot accommodate \mathbf{U} and \mathbf{V}^\dagger at the same time.) The type of `gesvd` can be given using the above definition of SVD flags:

```
val gesvd : jobu:('u_cols, 'm, ('m, 'n) min, ('m, 'n) min, z) svd_job →
          jobvt:('vt_cols, 'm, ('m, 'n) min, ('m, 'n) min, z) svd_job →
          ?s:(('m, 'n) min, cnt) vec      (* D *) →
          ?u:(('m, 'u_cols, 'u_cd) mat   (* U *) →
          ?vt:(('vt_rows, 'n, 'vt_cd) mat (* V† *) →
          ('m, 'n, 'a_cd) mat           (* A *) →
          (('m, 'n) min, 's_cd) vec      (* D *) *
          ('m, 'u_cols, 'u_cd) mat      (* U *) *
          ('vt_rows, 'n, 'vt_cd) mat    (* V† *)
```

In accordance with the type-level trick described in the next section, \mathbf{D} of `gesdd` and `gesvd` has type `(('m, 'n) min, cnt) vec` in the argument (contravariant position) and type `(('m, 'n) min, 's_cd) vec` in the return value (covariant position) because \mathbf{D} refers to a contiguous memory region.

3.2 Subtyping for discrete memory access

In Fortran, elements of a matrix are stored in column-major order in a flat, contiguous memory region. BLAS and LAPACK functions can take part of a matrix (like a row, a column, or a submatrix) and use it for computation without copying the elements, so they need to access the memory discretely in order to access the elements. However, some original functions of LACAML do not support such discrete access. For compatibility and soundness, we need to prevent those functions from receiving (sub)matrices that require discrete accesses while allowing the converse (i.e., the other functions may receive contiguous matrices as well as discrete ones). We achieved this by extending the type definition of matrices by adding a third parameter for “contiguous or discrete” flags (in addition to the existing two parameters for dimensions):

```
type ('m, 'n, 'cnt_or_dsc) mat (* 'm-by-'n matrices *)
type cnt (* phantom *)
type dsc (* phantom *)
```

Then, formal arguments that may be either contiguous or discrete matrices are given type `('m, 'n, 'cnt_or_dsc) mat`, while the types of (formal) arguments that *must* be contiguous are specified in the form `('m, 'n, cnt) mat`. In contrast, return values that may be either contiguous or discrete have type `('m, 'n, dsc) mat`, while those that are known to be always contiguous are typed `('m, 'n, 'cnt_or_dsc) mat` so that they can be mixed with discrete matrices.

Appendix A presents a generalization of the idea in this subsection to encode subtyping via phantom types.

3.3 Dynamic checks remaining

Although many high-level operations provided by BLAS and LAPACK can be statically verified by using the scheme described above as far as equalities of dimensions are concerned, other operations still require dynamic checks for inequalities:

- `get` and `set` operations allow accesses to an element of a matrix via given indices, which must be less than the dimensions of the matrix. This inequality is dynamically checked and no static

size constraint is imposed. (The use of these low-level functions is therefore unrecommended and high-level matrix operations such as `map` or BLAS/LAPACK functions should be used instead; see Section 4.2 for details.)

- As mentioned above, BLAS and LAPACK functions can take a submatrix without copying it. Our original function `submat` returns such a submatrix for the dimensions given. This submatrix must be smaller than the original matrix.
- An efficient memory arrangement for band matrices in BLAS and LAPACK, called the *band storage scheme*, requires a dynamic check that the specified numbers of sub- and super-diagonals are smaller than the size of the band matrix; for details, see Section 3.4.2.
- There are several high-level functions with specifications that essentially involve submatrices and inequalities of indices, such as `syevr` (for finding eigenvalues), `orgqr`, and `ormqr` (for QR factorization).
- For most LAPACK functions, the workspace for computation can be given as an argument. It must be larger than the minimum workspace required as determined by each function (and other arguments).

We gave these dynamically checked functions the suffix `_dyn`, like `get_dyn` and `set_dyn` (with the exception of the last bullet point since almost *all* LAPACK functions can take the workspace as a parameter whose size must be checked dynamically).

These inequalities are dynamically checked by our library interface for the sake of usability and compatibility with LACAML because the interface would become too complex if the inequalities were represented as types: for example, consider `('m, 'n) le` as a type that would represent $'m \leq 'n$; then the types of `get` and `set` could be given as follows (one is the type for the natural number 1):

```
val get : (one, 'i) le → ('i, 'n) le → ('n, _) vec → 'i size → float
val set : (one, 'i) le → ('i, 'n) le → ('n, _) vec → 'i size → float →
  unit
```

Not only the users must pass two extra arguments, but they would also have to *derive* the inequalities by applying functions for axioms such as reflexivity and transitivity. We have rejected this approach because, after all, users want to write linear algebra, not proof terms.

3.4 Inequality capabilities

Although most inequalities are dynamically checked, we introduced capabilities and types for guaranteeing inequalities in the following two cases as they arise naturally.

3.4.1 Submatrices and subvectors

All BLAS and LAPACK functions support operations on subvectors and submatrices. For instance,

```
lange ~m ~n ~ar ~ac a
```

computes the norm of the m -by- n submatrix in matrix `a` where element (i, j) corresponds to the $(i + ar - 1, j + ac - 1)$ element of `a`³. We cannot statically verify whether this function call is safe, i.e., the submatrix is indeed a submatrix of `a`. Since adding dynamic checks to *all* BLAS and LAPACK functions is undesirable, we have introduced a new, separate function

³The actual `lange` function takes another parameter `norm` that indicates the kind of the norm, e.g., 1-norm, Frobenius norm, etc.

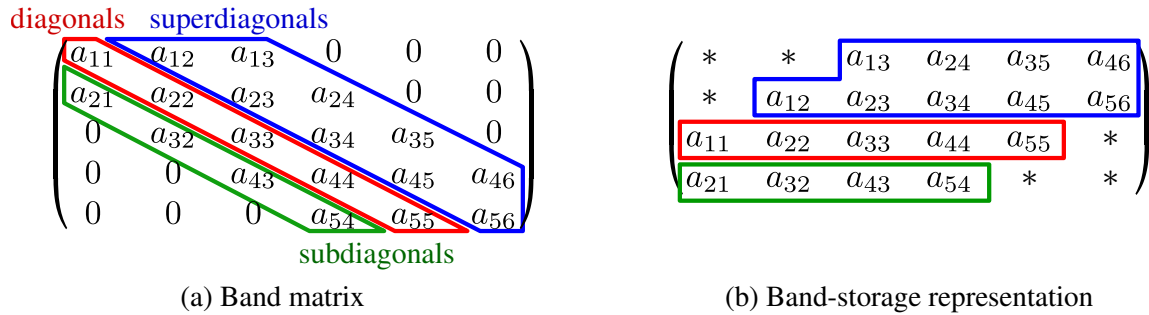


Figure 2: Band storage scheme

```

val submat_dyn : 'm size → 'n size → ?ar:int → ?ac:int →
    ('k, 'l, 'cnt_or_dsc) mat → ('m, 'n, dsc) mat

```

to return a submatrix of the given matrix. Then the type of `lange` can be given simply as:

```

val lange : ('m, 'n, 'cnt_or_dsc) mat → float

```

Thus `lange` itself requires no dynamic checks because the inequalities are already checked by `submat_dyn`.

Similarly, we defined the function `subvec_dyn` to return a subvector of a given vector or matrix.

3.4.2 Band storage scheme

In BLAS (and therefore LAPACK), an m -by- n band matrix with kl subdiagonals and ku superdiagonals can be stored in a matrix with $kl + ku + 1$ rows and n columns. This *band storage scheme* is used in practice only when $kl, ku \ll \min(m, n)$. The (i, j) element in the original matrix is stored in the $(ku + 1 + i - j, j)$ element of the band-storage representation, where $\max(1, j - ku) \leq i \leq \min(m, j + kl)$. Figure 2 shows a 5-by-6 band matrix with two superdiagonals and one subdiagonal, and its band-storage representation as an example (the `*` symbol denotes an unused element). A matrix like (b) is passed to special functions like `gbmv` (which multiplies a band matrix to a vector) for band-storage representations.

We implemented our original function `geband_dyn`⁴ that converts a band matrix into band-storage representation

```

val geband_dyn : 'kl size → 'ku size → ('m, 'n) mat →
    (('m, 'n, 'kl, 'ku) geband, 'n) mat

```

where the phantom type `('m, 'n, 'kl, 'ku) geband` guarantees the inequalities $kl < m$ and $ku < n$, and represents the height of the band-storage representation of an m -by- n band matrix with kl subdiagonals and ku superdiagonals. `geband_dyn` dynamically checks the inequalities and performs the conversion.

The matrix-vector multiplication function `gbmv` in LACAML⁵ is typed as

```

val gbm : ?m:int → ?beta:num_type →
    ?y:vec → (* y *)

```

⁴`geband_dyn` means GEneral BAND matrix. There is also a variant `syband_dyn` for symmetric band matrices, hence the name.

⁵We omit an optional argument `?n:int` to take a submatrix of `A`, since we separate such submatrix operations as in Section 3.4.1.

```
?trans:trans3 →
?alpha:num_type →
mat (* A *) →
int (* kl *) → int (* ku *) → vec (* x *) → vec (* y *)
```

while we give it the following type using `('m, 'n, 'kl, 'ku) geband`.

```
val gbmv : m:'a_m size → ?beta:num_type →
  ?y:'m vec (* y *) →
  trans:(('a_m, 'a_n) mat → ('m, 'n) mat) trans →
  ?alpha:num_type →
  (('a_m, 'a_n, 'kl, 'ku) geband, 'a_n) mat (* A *) →
  'kl size (* kl *) → 'ku size (* ku *) → 'n vec (* x *) → 'm vec (* y *)
```

This function computes $\mathbf{y} := \alpha\mathbf{A}\mathbf{x} + \beta\mathbf{y}$ or $\mathbf{y} := \alpha\mathbf{A}^\top\mathbf{x} + \beta\mathbf{y}$ for an m -by- n matrix \mathbf{A} , an n -dimensional vector \mathbf{x} , and an m -dimensional vector \mathbf{y} . Note that `gbmv` requires no dynamic check of $kl < m$ or $ku < n$ because the inequalities are statically guaranteed by the type `('m, 'n, 'kl, 'ku) geband`.

4 Porting of OCaml-GPR

To evaluate the usability of SLAP, we ported OCaml-GPR—a practical machine learning library for Gaussian process regression, written using LACAML by the same author—to use SLAP. The ported library, called SGPR, is available at <https://github.com/akabe/sgpr>.

Just to give a (very) rough feel for the library via a simple (but non-trivial) example, the type of a function that “calculates a covariance matrix of inputs given a kernel”

```
val calc_upper : Kernel.t → Inputs.t → mat
```

is augmented like

```
val calc_upper : ('D, _, _) Kernel.t → ('D, 'n) Inputs.t →
  ('n, 'n, 'cnt_or_dsc) mat
```

indicating that it takes $'n$ vectors of dimension $'D$ and returns an $'n$ -by- $'n$ contiguous matrix.

We added comments of the form `(#! label *)` on lines changed in the SGPR source code to investigate the categories and numbers of changes required for the porting, where *label* corresponds to a category of changes, e.g., `(#! IDX *)` for replacement of index-based accesses. We classified the changes into 19 categories, some of which we outline here (see <https://akabe.github.io/sgpr/changes.pdf> for the others).

4.1 Changes that could be made mechanically

Of the 19 categories of changes required, 12 could be made mechanically. Here we describe three representative examples.

Replacement of index-based accesses (IDX) In LACAML, the syntax sugar `x.{i, j}` can be used for index-based accesses to elements of vectors and matrices because they are implemented with the built-in OCaml module `Bigarray`. In SLAP, however, this syntax sugar cannot be used since `('n, 'cd) vec` and `('m, 'n, 'cd) mat` are abstract⁶; one must use the `get_dyn` and `set_dyn` functions instead.

⁶Since the development version of OCaml (4.03, not released yet) will support user-defined index operators `.{ }` and `.{ , }`, our replacement functions `get_dyn` and `set_dyn` will be unnecessary in the (near) future.

Rewriting of the flags (RF) The transpose flags had to be rewritten from `'N`, `'T`, and `'C` to `normal`, `trans`, and `conjtr` (and similarly for side and SVD job flags) for the sake of typing as described in Section 3.1.

Insertion of type parameters (ITP) Recall that we changed the types `vec` and `mat` on the right hand side of a type definition to `('n, 'cd) vec` and `('m, 'n, 'cd) mat`, respectively. It is then necessary to add the type parameters `'m`, `'n`, and `'cd` on the left hand side as well. Theoretically, it suffices to give fresh type parameters to all `vec` and `mat`. For example⁷,

```
module M : sig
  type t
  val f : int → t
end = struct
  type t = {
    n : int;
    id : mat;
  }
  let f n = { n; id = Mat.identity n }
end
```

should be rewritten:

```
module M : sig
  type ('a, 'b, 'c, 'd) t
  val f : 'a size → ('a, 'a, 'a, 'cnt_or_dsc) t
end = struct
  type ('a, 'b, 'c, 'd) t = {
    n : 'a size;
    id : ('b, 'c, 'd) mat;
  }
  let f n = { n; id = Mat.identity n }
end
```

Note that, in the latter code, type parameters of function `f` could be automatically inferred by OCaml.

In practice, however, it would introduce too many type parameters to take *all* of them fresh. We thus unified the type parameters that are known to be always equal by looking at the constructor functions such as `f`:

```
module M : sig
  type ('n, 'cnt_or_dsc) t (*! ITP *)
  val f : 'n size → ('n, 'cnt_or_dsc) t (*! ITP *)
end = struct
  type ('n, 'cnt_or_dsc) t = { (*! ITP *)
    n : 'n size; (*! ITP *)
    id : ('n, 'n, 'cnt_or_dsc) mat; (*! ITP *)
  }
  let f n = { n; id = Mat.identity n }
end
```

⁷This example is imaginary and *not* from the real OCaml-GPR code, which is too complex to explain on paper.

4.2 Changes that had to be made manually

Other changes needed human brain and had to be made manually. We here explain the following two representatives of the seven categories of manual changes. (The other categories are function types that depend on the values of arguments, and some ad hoc changes; again see <https://akabe.github.io/sgpr/changes.pdf>.)

Insertion of type annotations (ITA) When a matrix operation is implemented by using low-level index-based access functions, its size constraints cannot be inferred statically (since they are checked only at runtime). For example, consider the function `axby`, which calculates $\alpha\mathbf{x} + \beta\mathbf{y}$ where α and β are scalar values, and \mathbf{x} and \mathbf{y} are vectors.

```
let axby alpha x beta y =
  let n = Vec.dim x in
  let z = Vec.create n in
  for i = 1 to Size.to_int n do
    Vec.set_dyn z i
      (alpha *. (Vec.get_dyn x i) +. beta *. (Vec.get_dyn y i))
  done;
  z
```

The dimensions of vectors \mathbf{x} and \mathbf{y} must be the same, but OCaml cannot infer that:

```
val axby : float → ('n, _) vec → float → ('m, _) vec → ('n, _) vec
```

There are two ways to solve this problem. One is to type-annotate the function by hand:

```
let axby alpha (x : ('n, _) vec) beta (y : ('n, _) vec) =
  ...
```

The other is to use high-level operations such as BLAS/LAPACK functions or `map2` instead of the low-level operations `get_dyn` and `set_dyn`:

```
let axby alpha x beta y =
  let z = copy y in (* z := y *)
  scal beta z;      (* z := beta * z *)
  axpy ~alpha ~x z; (* z := alpha * x + z *)
  z
```

or

```
let axby alpha x beta y =
  Vec.map2 (fun xi yi → alpha *. xi +. beta *. yi) x y
```

In either way, we need to rewrite existing programs by considering their meanings.

We encountered five such functions in OCaml-GPR and adopted the former approach for all of them so as to keep the changes minimal.

Escaping generative phantom types (EGPT) We needed to prevent a generative phantom type from escaping its scope. Consider the following function implemented using LACAML⁸.

```
let vec_of_array a = Vec.init (Array.length a) (fun i → a.(i))
```

⁸Again, this is just an example for pedagogy; SLAP itself supplies a similar function `Vec.of_array`.

It converts an array into a vector. In SLAP, the above code causes a type error because `Vec.init` expects a size value of singleton type `'n size`, not an integer, as the first argument. We thus should convert the integer into a size with `Size.of_int_dyn : int → (module SIZE)9` as follows:

```

module type SIZE = sig
  type n
  val value : n size
end

let vec_of_array a =
  let module N = (val Size.of_int_dyn (Array.length a) : SIZE) in
  Vec.init N.value (fun i → a.(i))

```

The fix may seem correct, but it does not type-check in OCaml because the generative phantom type `N.n` escapes its scope.

There are two ways to handle this situation in SLAP. One is to insert the argument `n` for the size of the array and remove the generative phantom type from the function:

```

let vec_of_array n a = (* : 'n size → float array → ('n, _) vec *)
  if Size.to_int n <> Array.length a then invalid_arg "error";
  Vec.init n (fun i → a.(i))

```

In this approach, the generative phantom type should be given from the outside of the function. For example, `vec_of_array` can be called like:

```

let f a =
  let module N = (val Size.of_int_dyn (Array.length a) : SIZE) in
  let v = vec_of_array N.value a in
  printf "%a" pp_rfvec v

```

However, the dynamic check in the definition of `vec_of_array` is redundant in this case.

The other way is to define `vec_of_array : float array → (?, _) vec` using a first-class module (or a functor):

```

let vec_of_array a = (* : float array → (module VEC) *)
  let module N = (val Size.of_int_dyn (Array.length a) : SIZE) in
  (module struct
    type n = N.n
    let value = Vec.init N.value (fun i → a.(i))
  end : VEC)

```

In this approach, the generative phantom type is created inside the function. No extra dynamic check is required, but the user needs to write a type annotation for the returned module like `(val vec_of_array a : VEC)` at the caller site. We suppose that conversion from the first `vec_of_array` implemented in LACAML into the last code can be made automatically by inserting packing and unpacking of first-class modules: packing should be inserted where a generative phantom type escapes, and unpacking should be inserted where the contents (a vector or a matrix) of the packed module are used. It is however burdensome to manually insert them because of the heavy syntax. We thus adopted the former approach for our (manual) porting.

Note that, in either approach, the conversion may introduce another escape of the generative phantom type at the caller site and therefore may have to be repeated (until it reaches the main routine in the worst case, though we conjecture from our experiences with SLAP that such cases are rare).

⁹ It raises an exception if the given integer is negative because a size must be non-negative, hence the suffix `_dyn`.

Table 1: Number and percentage of mechanically changed lines

	S2I	SC	SOP	I2S	IDX	RF	IF	SUB	ETA	RID	RMDC	ITP	Total
lib/block_diag.mli	0	0	0	0	0	0	0	0	0	1	0	5	6
lib/block_diag.ml	1	0	0	0	0	0	0	0	0	1	6	1	9
lib/cov_const.mli	0	0	0	0	0	0	0	0	0	0	0	5	5
lib/cov_const.ml	2	0	0	0	1	0	0	0	0	2	0	9	14
lib/cov_lin_one.mli	0	0	0	0	0	0	0	0	0	1	0	5	6
lib/cov_lin_one.ml	0	0	0	0	1	4	2	0	0	2	0	9	17
lib/cov_lin_ard.mli	0	0	0	0	0	0	0	0	0	1	0	5	6
lib/cov_lin_ard.ml	7	0	0	0	10	5	2	0	0	2	0	9	32
lib/cov_se_iso.mli	0	0	0	0	0	0	0	0	0	1	0	5	6
lib/cov_se_iso.ml	18	4	0	0	31	0	0	0	0	2	2	14	71
lib/cov_se_fat.mli	0	0	0	0	0	0	0	0	0	1	0	10	11
lib/cov_se_fat.ml	43	9	1	0	87	2	2	0	0	2	8	23	174
lib/fitc_gp.mli	0	0	0	0	0	0	0	0	0	0	0	0	0
lib/fitc_gp.ml	81	3	3	0	63	19	26	14	34	3	15	69	298
lib/interfaces.ml	0	0	0	0	0	0	0	0	0	1	0	196	197
lib/gpr_utils.ml	10	0	0	0	13	0	2	0	0	4	17	1	46
app/ocaml_gpr.ml	13	0	2	4	10	0	0	0	0	3	0	8	35
Total	175	16	6	4	216	30	34	14	34	27	48	374	933
Percentage	2.89	0.26	0.10	0.07	3.56	0.49	0.56	0.23	0.56	0.45	0.79	6.17	15.39

4.3 Results

Table 1 shows the number of lines that required mechanical changes. The major change was ITP (6.17 %) because OCaml-GPR consists of several large modules with a large number of functions involving sized types. Most of the ITP changes have been made in `lib/interfaces.ml`, which defines all the signatures provided by OCaml-GPR. IDX was the second largest (3.56 %) because index-based access functions are frequently used in OCaml-GPR, even when they could be replaced with high-level matrix operations such as `map`.

Table 2 shows the numbers and percentages of lines for which the required changes had to be made manually, and Table 3 gives the total of all changes. Overall, 18.4 % of lines required some changes, out of which (with some overlap) 15.4 % were mechanical and 3.6 % required human brain. From these results, we conjecture in general that the number of non-trivial changes required for a user program of SLAP is small.

5 Related work

Dependent ML (DML) [22] and sized type [3] can statically verify consistency among the sizes of collections such as lists and arrays, using dependent types on natural numbers. In ATS [4], a successor of DML, BLAS and LAPACK bindings are provided. Advantages of dependent types over our approach are: (1) they can represent more complex specifications including inequalities such as array bounds; (2) they can verify the consistency of sizes in the *internal* implementations of vector and matrix operations (though the BLAS and LAPACK bindings for ATS are currently implemented as wrappers of C functions, so the internals are not statically verified). Conversely, our approach only requires fairly standard ML types and module system, and application programs can be ported almost mechanically, while dependent types generally require non-trivial changes to the programming language and application programs.

The dimensions of vectors and matrices can also be represented [11] using GADT, a lightweight

Table 2: Number and percentage of manually changed lines

	ITA	EGPT	O2L	FT	ET	DKS	FS	Total
lib/block_diag.mli	0	0	0	0	0	0	0	0
lib/block_diag.ml	0	0	0	0	0	0	0	0
lib/cov_const.mli	0	0	0	0	0	0	2	2
lib/cov_const.ml	0	1	0	0	0	1	8	10
lib/cov_lin_one.mli	0	0	0	0	0	0	2	2
lib/cov_lin_one.ml	0	0	0	0	0	1	12	13
lib/cov_lin_ard.mli	0	0	0	0	0	0	2	2
lib/cov_lin_ard.ml	0	0	0	0	0	1	8	9
lib/cov_se_iso.mli	0	0	0	0	0	0	2	2
lib/cov_se_iso.ml	2	0	0	0	0	1	6	9
lib/cov_se_fat.mli	0	0	0	4	0	0	0	4
lib/cov_se_fat.ml	0	0	0	28	0	1	1	30
lib/fitc_gp.mli	0	0	0	0	0	0	0	0
lib/fitc_gp.ml	0	28	31	16	0	0	0	68
lib/interfaces.ml	0	11	7	5	0	3	0	26
lib/gpr_utils.ml	6	0	0	0	0	0	1	7
app/ocaml_gpr.ml	0	17	0	6	16	0	0	35
Total	8	57	38	59	16	8	44	219
Percentage	0.13	0.94	0.63	0.97	0.26	0.13	0.73	3.61

Table 3: Number and percentage of all changed lines

	Lines	Mechanical	Manual	Total
lib/block_diag.mli	56	6	0	6
lib/block_diag.ml	58	9	0	9
lib/cov_const.mli	52	5	2	6
lib/cov_const.ml	141	14	10	16
lib/cov_lin_one.mli	56	6	2	7
lib/cov_lin_one.ml	149	17	13	26
lib/cov_lin_ard.mli	56	6	2	7
lib/cov_lin_ard.ml	188	32	9	39
lib/cov_se_iso.mli	58	6	2	7
lib/cov_se_iso.ml	343	71	9	78
lib/cov_se_fat.mli	105	11	4	15
lib/cov_se_fat.ml	680	174	30	199
lib/fitc_gp.mli	151	0	0	0
lib/fitc_gp.ml	2294	298	68	364
lib/interfaces.ml	1008	197	26	215
lib/gpr_utils.ml	229	46	7	53
app/ocaml_gpr.ml	440	35	35	66
Total	6064	933	219	1113
Percentage	100.00	15.39	3.61	18.35

form of dependent types. Existential types can be implemented not only using first-class modules but also using GADT.

The idea of using phantom and generative types for static size checking is not novel. Kiselyov and Shan [12] implemented DML-like size checking (including inequalities, e.g., array bound checking) by CPS encoding of existential types using first-class polymorphism. Eaton [6] developed a linear algebra library with static size checking for matrix operations as a “strongly statically typed” binding of GSL-Haskell¹⁰. His basic idea is similar to ours, but he adopted Template Haskell [19] for the CPS encoding of generative types. The approaches of [6, 12] need CPS conversion when a generative type escapes its scope and thereby change the structures of programs. In contrast, we either implemented generative types with first-class modules in OCaml instead of the CPS encoding, or else removed them in the first place, like the conversions in EGPT. Our contribution is the discovery that practical size checking for a linear algebra library can be constructed on the simple idea of verifying mostly the equality of sizes without significantly restructuring application programs.

Braibant and Pous [2] implemented static size checking of matrix operations using phantom types on Coq. It requires more type annotations than our interface.

Eigen [7] is another practical linear algebra libraries with static size checking. It does not statically check the consistency of dynamically determined sizes of matrices and vectors.

6 Conclusions

Our proposed linear algebra library interface SLAP uses generative phantom types to statically ensure that most operations on matrices satisfy dimensional constraints. It is based on a simple idea—only the equality of sizes needs to be verified—and can be realized by using a fairly standard type and module system of ML. We implemented this interface on top of LACAML and then ported OCaml-GPR to it. Most of the high-level matrix operations in the BLAS and LAPACK linear algebra libraries were successfully typed, and few non-trivial changes were required for the porting.

We did not find any bug in LACAML or OCaml-GPR, maybe because both libraries have already been well tested and debugged or carefully written in the first place. However, in our experience of implementing other (relatively small) programs¹¹, our version of the libraries have been particularly useful when developing a new library or application on top since they detect an error not only *earlier* (i.e., at compile time instead of runtime) but also *at higher level*: for instance, if the programmer misuses a function of SGPR, an error is reported at the caller site rather than somewhere deep inside the call stack from the function.

Interesting directions for future work include formalization of the idea of generative phantom types, and extension of the static types to enable verification of inequalities (in addition to equalities), just to name a few.

References

- [1] Matthias Blume (2001): *No-Longer-Foreign: Teaching an ML compiler to speak C “natively”*. *Electr. Notes Theor. Comput. Sci.* 59(1), pp. 36–52. Available at [http://dx.doi.org/10.1016/S1571-0661\(05\)80452-9](http://dx.doi.org/10.1016/S1571-0661(05)80452-9).

¹⁰ GSLHaskell is a binding of the GNU Scientific Library (GSL) [9], a library for linear algebra and numerical computation on C and C++. GSL also provides interfaces for BLAS and LAPACK.

¹¹such as neural networks; see <https://github.com/akabe/slap/tree/master/examples>

- [2] Thomas Braibant & Damien Pous (2010): *An Efficient Coq Tactic for Deciding Kleene Algebras*. In: *Interactive Theorem Proving, First International Conference, ITP 2010, Edinburgh, UK, July 11-14, 2010. Proceedings, Lecture Notes in Computer Science 6172*, Springer, pp. 163–178. Available at http://dx.doi.org/10.1007/978-3-642-14052-5_13.
- [3] Wei-Ngan Chin & Siau-Cheng Khoo (2001): *Calculating Sized Types*. *Higher-Order and Symbolic Computation* 14(2-3), pp. 261–300. Available at <http://dx.doi.org/10.1023/A:1012996816178>.
- [4] Sa Cui, Kevin Donnelly & Hongwei Xi (2005): *ATS: A Language That Combines Programming with Theorem Proving*. In: *FroCoS, Lecture Notes in Computer Science 3717*, Springer, pp. 310–320. Available at http://dx.doi.org/10.1007/11559306_19.
- [5] Olivier Danvy (1998): *Functional Unparsing*. *J. Funct. Program.* 8(6), pp. 621–625. Available at <http://dx.doi.org/10.1017/S0956796898003104>.
- [6] Frederik Eaton (2006): *Statically typed linear algebra in Haskell*. In: *Proceedings of the ACM SIGPLAN Workshop on Haskell, Haskell 2006, Portland, Oregon, USA, September 17, 2006*, ACM, pp. 120–121. Available at <http://dx.doi.org/10.1145/1159842.1159859>.
- [7] *Eigen*. <http://eigen.tuxfamily.org/>.
- [8] Matthew Fluet & Riccardo Pucella (2006): *Phantom types and subtyping*. *J. Funct. Program.* 16(6), pp. 751–791. Available at <http://dx.doi.org/10.1017/S0956796806006046>.
- [9] Mark Galassi et al.: *the GNU Scientific Library (GSL)*. <http://www.gnu.org/software/gsl/>.
- [10] Ralf Hinze (2003): *Fun with phantom types*. In Jeremy Gibbons & Oege de Moor, editors: *The Fun of Programming*, Cornerstones of Computing, Palgrave Macmillan, pp. 245–262.
- [11] hyone: *Length Indexed Matrix and Indexed Functor*. <https://gist.github.com/hyone/3990929>.
- [12] Oleg Kiselyov & Chung chieh Shan (2007): *Lightweight Static Capabilities*. *Electr. Notes Theor. Comput. Sci.* 174(7), pp. 79–104. Available at <http://dx.doi.org/10.1016/j.entcs.2006.10.039>.
- [13] Daan Leijen & Erik Meijer (1999): *Domain specific embedded compilers*. In: *Proceedings of the Second Conference on Domain-Specific Languages (DSL '99), Austin, Texas, USA, October 3-5, 1999*, ACM, pp. 109–122. Available at <http://dx.doi.org/10.1145/331960.331977>.
- [14] Markus Mottl: *OCaml-GPR – Efficient Gaussian Process Regression in OCaml*. <https://github.com/mmottl/gpr>.
- [15] Markus Mottl & Christophe Troestler: *LACAML – Linear Algebra for OCaml*. <https://github.com/mmottl/lacaml>.
- [16] NetLib: *BLAS (Basic Linear Algebra Subprograms)*. <http://www.netlib.org/blas/>.
- [17] NetLib: *LAPACK – Linear Algebra PACKage*. <http://www.netlib.org/lapack/>.
- [18] Stanford University’s Pervasive Parallelism Laboratory (PPL): *OptiML*. <http://stanford-ppl.github.io/Delite/optiml/>.
- [19] Tim Sheard & Simon L. Peyton Jones (2002): *Template meta-programming for Haskell*. *SIGPLAN Notices* 37(12), pp. 60–75. Available at <http://dx.doi.org/10.1145/636517.636528>.
- [20] Arvind K. Sujeeth, HyoukJoong Lee, Kevin J. Brown, Tiark Rompf, Hassan Chafi, Michael Wu, Anand R. Atreya, Martin Odersky & Kunle Olukotun (2011): *OptiML: An Implicitly Parallel Domain-Specific Language for Machine Learning*. In: *Proceedings of the 28th International Conference on Machine Learning, ICML 2011, Bellevue, Washington, USA, June 28 - July 2, 2011*, Omnipress, pp. 609–616.
- [21] *uBlas*. http://www.boost.org/doc/libs/1_55_0/libs/numeric/ublas/doc/.
- [22] Hongwei Xi (2007): *Dependent ML – An approach to practical programming with dependent types*. *J. Funct. Program.* 17(2), pp. 215–286. Available at <http://dx.doi.org/10.1017/S0956796806006216>.

Table 4: Encoding of supertype \top and subtype \perp

	Positive (covariant) position		Negative (contravariant) position
Supertype \top	\perp	\top	$\top \text{_or_} \perp$
Subtype \perp	$\top \text{_or_} \perp$	\perp	\perp

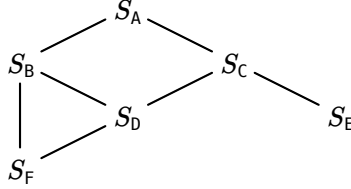


Figure 3: Example of powerset lattice

A Encoding of subtyping

We explain generalization of the subtyping encoding used in Section 3.2. To start with, consider a subtyping with only one base case $\top \text{ :> } \perp$. Table 4 shows encoding of the types \top and \perp (which depends on their positions of appearance). It is similar to the subtyping encoding for the types of contiguous and discrete matrices (without type parameters \top_m and \perp_n for dimensions) in Section 3.2. \perp and \perp are phantom types, and $\top \text{_or_} \perp$ is the sum type of \top and \perp where $\top \text{_or_} \perp$ is a phantom type parameter.

More generally, to encode an arbitrary subtyping hierarchy with a finite number of base cases, we give an encoding of *powerset lattices*. The powerset lattice of a finite set S is the lattice of all subsets of S , ordered by inclusion. As in [8], we represent our subtyping relation as the inclusion relation between subsets of some S so that $\top \text{ :> } \perp$ iff $S_\top \supseteq S_\perp$ for any types \top and \perp , where S_\top and S_\perp are some appropriate subsets of S corresponding to \top and \perp , respectively. For example, consider the subtyping relation illustrated in Figure 3, where A is the largest type, and E and F are the smallest. Let S be $\{1, 2, 3, 4\}$. Then we can take, e.g., $S_A = \{1, 2, 3, 4\}$, $S_B = \{1, 2\}$, $S_C = \{1, 3, 4\}$, $S_D = \{2, 3\}$, $S_E = \{4\}$, and $S_F = \{1\}$.

We now assume a total ordering s_1, s_2, \dots, s_n among the elements of S , where n is the cardinality of S . A base type \top is encoded as $\langle \top \rangle_+ \tau$ at covariant positions, and as $\langle \top \rangle_- \tau$ at contravariant positions, where $\langle \top \rangle_+$ and $\langle \top \rangle_-$ are n -tuple types defined as

$$\langle \top \rangle_+ \stackrel{\text{def}}{=} t_1 * \dots * t_n \quad \text{where } t_i = \begin{cases} \top & \text{if } s_i \in S_\top \\ \top \text{_or_} \perp & \text{otherwise,} \end{cases}$$

$$\langle \top \rangle_- \stackrel{\text{def}}{=} t_1 * \dots * t_n \quad \text{where } t_i = \begin{cases} \top \text{_or_} \perp & \text{if } s_i \in S_\top \\ \perp & \text{otherwise,} \end{cases}$$

for phantom types \top and \perp . We require that every type parameter $\top \text{_or_} \perp$ is fresh in *each* $\langle \cdot \rangle_+$ and $\langle \cdot \rangle_-$.

Table 5 shows the encoding of the subtyping relation in Figure 3. We can verify, e.g., that $\langle E \rangle_+$ can be unified with $\langle A \rangle_-$, $\langle C \rangle_-$, or $\langle E \rangle_-$, but not with $\langle B \rangle_-$, $\langle D \rangle_-$, or $\langle F \rangle_-$. That is, a value of type E can be passed to a function as an argument of type A , C , or E , but not as B , D , or F . In addition, if we put values of type E and F in the same list, it can be passed as an argument of type A *list* or C *list* (i.e., a list of elements with a common supertype of E and F).

We implement Bot , a subtype of any type, as a single type parameter (i.e., $\forall \alpha. \alpha$), which can

Table 5: Encoding of the subtyping hierarchy of Figure 3

	$\langle \cdot \rangle_+$				$\langle \cdot \rangle_-$			
A	w	*	w	*	w	*	w	' a ₁ * ' a ₂ * ' a ₃ * ' a ₄
B	w	*	w	*	' a ₃ * ' a ₄	' a ₁ * ' a ₂ * z * z		
C	w	*	' a ₂ * w * w	' a ₁ * z * ' a ₃ * ' a ₄				
D	w	*	' a ₂ * w * ' a ₄	' a ₁ * z * ' a ₃ * z				
E	' a ₁ * ' a ₂ * ' a ₃ * w	z * z * z * ' a ₄						
F	w	*	' a ₂ * ' a ₃ * ' a ₄	' a ₁ * z * z * z				

be instantiated to any type. By replacing $' a_i$ with Bot in the definitions of $\langle \cdot \rangle_+$ and $\langle \cdot \rangle_-$, we obtain $U <: T \iff \langle U \rangle_+ <: \langle T \rangle_+ \wedge \langle U \rangle_- >: \langle T \rangle_-$.

Related work Fluet and Pucella [8] also proposed a subtyping scheme using phantom types on the ML type system. They focused on an encoding of Hindley-Milner polymorphism extended with subtyping. Their approach can encode a type like $\forall \alpha <: T. \alpha \rightarrow \alpha$, but does not achieve the contravariance of argument types. In contrast, our approach accomplishes both covariance and contravariance while it does not support universal types.

B Generalization of flag-dependent function types

We generalize the phantom type trick in Section 3.1 for function types that depend on flags (cf. Danvy's typing of `printf` [5]). The type of the following function depends on values of x_1, \dots, x_n ($n = 1$ in many functions of BLAS and LAPACK, but several functions such as `gemm` and `gesvd` take two or more flags).

```
val f :  $\Pi x_1:t_1. \Pi x_2:t_2. \dots \Pi x_n:t_n. (\mathcal{T}_1(x_1), \mathcal{T}_2(x_2), \dots, \mathcal{T}_n(x_n)) \text{ u}$ 
```

\mathcal{T}_i is a function that maps a flag value to an ML type, $(\alpha_1, \dots, \alpha_n) \text{ u}$ is an ML type, and t_i is the type of the i th flags, e.g., [``N | `T | `C`] (for transpose flags), [``L | `R`] (for side flags) or [``A | `S | `O | `N`] (for SVD job flags). We assume that \mathcal{T}_i does not contain dependent types.

We consider how to type `f` without dependent types. First, we represent each type t_i and each flag v_{ij} of type t_i as follows:

```
type  $\alpha$  tt $i$  (* = t $i$  *)
val tt $i$ -v $ij$  :  $\mathcal{T}_i(v_{ij})$  tt $i$  (* = v $ij$  *)
```

Then `f` takes the flag representations $tt_{i-v_{ij}}$ as arguments of types α_i tt _{i} , thereby receiving the types $\mathcal{T}_i(v_{ij})$ tt _{i} as α_i :

```
val f :  $\forall \alpha_1, \alpha_2, \dots, \alpha_n. \alpha_1$  tt1  $\rightarrow \alpha_2$  tt2  $\rightarrow \dots \rightarrow \alpha_n$  tt $n$   $\rightarrow (\alpha_1, \alpha_2, \dots, \alpha_n) \text{ u}$ 
```

We show the encoding of the type of `gemm` as an example. The original type of it is

```
val gemm :  $\Pi \text{transa}:[\text{'N} | \text{'T} | \text{'C}] \rightarrow \Pi \text{transb}:[\text{'N} | \text{'T} | \text{'C}] \rightarrow$ 
  ( $\mathcal{T}(\text{transa}), \mathcal{T}(\text{transb})$ ) u
```

where

```
(('am, 'ak) mat  $\rightarrow$  ('m, 'k) mat, ('bk, 'bn) mat  $\rightarrow$  ('k, 'n) mat) u
= ?beta:num.type  $\rightarrow$  ?c:(('m, 'n) mat (* C *)  $\rightarrow$ 
```

```
?alpha:num_type → ('am, 'ak) mat (* A *) →
('bk, 'bn) mat (* B *) → ('m, 'n) mat (* C *)
```

and

$$\mathcal{T}(\text{trans}) = \begin{cases} ('m, 'n) \text{ mat} \rightarrow ('m, 'n) \text{ mat} & (\text{trans} = \text{'N}) \\ ('m, 'n) \text{ mat} \rightarrow ('n, 'm) \text{ mat} & (\text{trans} = \text{'T}, \text{'C}). \end{cases}$$

Our representations of the flags are

```
type  $\alpha$  trans (* = ['N | 'T | 'C] *)
val normal :  $\mathcal{T}(\text{'N})$  trans (* = 'N *)
val trans :  $\mathcal{T}(\text{'T})$  trans (* = 'T *)
val conjtr :  $\mathcal{T}(\text{'C})$  trans (* = 'C *)
```

and our type of gemm is:

```
val gemm :  $\forall \alpha_1, \alpha_2. \alpha_1$  trans  $\rightarrow \alpha_2$  trans  $\rightarrow (\alpha_1, \alpha_2)$  u
```

Side flags and SVD job flags can be represented similarly.