

# Data Parallel Idealised Algol

ROBERT ATKEY, University of Strathclyde  
MICHEL STEUWER, University of Glasgow  
SAM LINDLEY, The University of Edinburgh  
CHRISTOPHE DUBACH, The University of Edinburgh

---

Parallel hardware has the potential for accelerating a wide class of algorithms. However, expert programming skills are required to achieve maximum performance. Typical parallel architectures expose low-level hardware details through imperative programming interfaces where programmers explicitly encode device-specific optimisation strategies. This inevitably results in non-performance-portable programs that deliver suboptimal performance on other devices.

Functional programming models have recently seen a renaissance in the systems community as they offer solutions for tackling the performance-portability challenge. Recent work has demonstrated how to automatically choose high-performance parallelisation strategies for a wide range of hardware architectures encoded in a functional representation. However, the translation of such functional representations to the imperative program demanded by the hardware interface is typically performed ad hoc with no correctness guarantees and no guarantee to preserve the intended parallelisation strategy.

We introduce Data Parallel Idealised Algol (DPIA), a parallel dialect of Reynolds’ Idealised Algol, which supports both high-level functional code and low-level imperative code. We define and prove correct a translation from parallel high-level functional code to low-level data-race-free imperative code. Our performance results show that the translation yields performance on a par with hand-written low-level code.

CCS Concepts: • **Software and its engineering** → *Parallel programming languages; Functional languages; Imperative languages; Multiparadigm languages;*

## ACM Reference format:

Robert Atkey, Michel Steuwer, Sam Lindley, and Christophe Dubach. 2018. Data Parallel Idealised Algol. *Proc. ACM Program. Lang.* 0, 0, Article 0 (October 2018), 27 pages.  
<https://doi.org/Unpublished>

---

## 1 INTRODUCTION

Modern parallel hardware promises unprecedented performance, but only for those who can program it correctly. A single program must simultaneously describe the functional process for generating the result, as well as how to exploit the hardware’s resources to deliver it. These programs are commonly written in imperative languages like OpenMP, OpenCL, and CUDA, which are difficult to reason about, and difficult to refactor to try different parallelisation strategies.

To mitigate this problem, there has been an old, but now accelerating, trend towards functional languages for parallel computation. The programmer specifies the program they want in a high level language, using higher order functions and compositional abstractions. The compiler then takes responsibility for generating efficient low level imperative code. NESL [Blelloch 1993], Copper-Head [Catanzaro et al. 2011], LiquidMetal [Dubach et al. 2012], Accelerate [McDonnell et al. 2013], and Delite [Sujeeth et al. 2014] are examples. Since the “Sufficiently Smart Compiler” that always produces the highest performance code is yet to be fully realised, many modern systems have ways of describing the parallelisation strategy to be used either alongside or within the functional code. Such a strategy describes matters such as when and where parallelism ought to be used, and how

---

2018. 2475-1421/2018/10-ART0 \$X  
<https://doi.org/Unpublished>

50 data ought to be chunked for parallel processing. Ideally, we separate the concerns of functional  
51 behaviour from strategy and reason about each separately. Halide [Ragan-Kelley et al. 2018] is an  
52 example of a system where the strategy is given “on the side”. Lift [Steuwer et al. 2015, 2017] and  
53 parallel Haskell [Trinder et al. 1998] are examples of systems where strategies are embedded within  
54 the functional code. The latter approach has the advantage that we cannot write strategies that do  
55 not make sense for the computation that we wish to perform. We will primarily concern ourselves  
56 with this embedded strategy approach in this paper.

57 Given such a strategic functional programming system, the question is how to compile the  
58 combination down to the imperative code required by the hardware, preserving the given strategy,  
59 and maintaining correctness. In previous works, this question has been treated as a mere matter of  
60 code generation from the high level functional code. However, code generation can be subtle. It is  
61 easy to formulate code generation procedures that do not work for all inputs, or produce code that  
62 does not compile, or has race conditions. Moreover, it would be beneficial to know that the code  
63 generation process not only produces functionally correct code, but also strategically correct code.

64 In this paper, we formalise a *semantics and strategy preserving translation* from functional to  
65 parallel imperative code using *Data Parallel Idealised Algol* (DPIA), a novel variant of Reynolds’  
66 *Idealised Algol* [Reynolds 1997]. Idealised Algol orthogonally combines typed functional programming  
67 with imperative programming, a property that DPIA inherits. This allows us to start in the purely  
68 functional subset with a high level specification of the computation we wish to perform, and then  
69 to systematically rewrite it to “purely imperative” code (Section 4.1) that has a straightforward  
70 translation into OpenMP annotated C code that can be executed in parallel (Section 4.3). DPIA  
71 incorporates a substructural type system, incorporating ideas from Reynolds’ *Syntactic Control of*  
72 *Inteference* [O’Hearn et al. 1999; Reynolds 1978], which ensures that the generated programs are  
73 data race free. We show that our translation from functional to imperative is correct (Section 5).

74 Our use of Idealised Algol provides a greater benefit than just a correctness proof. Reynolds’  
75 Syntactic Control of Interference discipline provides a framework in which to explore the boundaries  
76 between functional and imperative parallel programming. We demonstrate the generality of our  
77 approach in Section 6 where we augment the basic design of DPIA with vectorisation and double  
78 buffering support. Our primary conceptual contribution is to promote Idealised Algol and its  
79 variants as expressive and flexible intermediate languages for compilation to low level imperative  
80 code that efficiently and correctly use the facilities provided by the hardware. In particular, syntactic  
81 interference control allows us to given types to the primitive hardware operations that stops us  
82 from incorrectly applying them. We explain the process of designing new primitives for DPIA in  
83 Section 3, when we introduce the type system.

84  
85 *Contributions.* This paper makes three major contributions:

- 86 (1) We introduce *Data Parallel Idealised Algol* (DPIA). This is an extension of O’Hearn et al.’s  
87 [1999] *Syntactic Control of Inteference Revisited* with indexed types, array and tuple data types,  
88 and primitives for data parallel programming. These extensions are fundamental for data  
89 parallel imperative programs. We describe DPIA in Section 3.
- 90 (2) We describe a formal translation from high level functional code annotated with parallelism  
91 strategies to parallel imperative code. We describe this translation fully in Section 4 and  
92 prove it functionally correct in Section 5.
- 93 (3) We have implemented our framework using OpenMP as the backend and evaluate the  
94 performed achieved by our generated code in Section 7.
- 95 (4) To demonstrate the flexibility of the DPIA framework, we extend our whole approach in Sec-  
96 tion 6 to include vectorisation and iterative computations implemented by double buffering.  
97

In the next section we introduce and motivate our approach and the use of Idealised Algol as a foundation. After the technical body of the paper in Sections 3-7, we discuss related work in Section 8 and our conclusions and future work in Section 9.

## 2 STRATEGY PRESERVING COMPILATION

We describe a compilation method that takes high level functional array code to low level parallel imperative code. Our approach is characterised by (a) expression of parallelisation strategies in high level functional code, where their semantic equivalence to the original code can be easily checked, and new strategies readily derived; and (b) a predictable and principled translation to low level imperative code that preserves parallelisation strategies.

### 2.1 Expressing Parallelisation Strategies in Functional Code

Here is an expression that describes the dot product of two vectors  $xs$  and  $ys$ :

$$\text{reduce } (+) \ 0 \ (\text{map } (\lambda x. \text{fst } x * \text{snd } x) \ (\text{zip } xs \ ys)) \quad (1)$$

This expression is a declarative specification of the dot product which does not specify how the computation should be performed. There are many possible strategies to execute the dot product computation in parallel. A naive strategy is to first perform the pairwise multiplication of  $xs$  and  $ys$  in parallel before sequentially summing up the intermediate results. We rewrite our expression to make this “parallel map, sequential reduce” strategy explicit:

$$\text{reduceSeq } (+) \ 0 \ (\text{mapPar } (\lambda x. \text{fst } x * \text{snd } x) \ (\text{zip } xs \ ys)) \quad (2)$$

Such a naive strategy is not always best. If we try to execute one parallel job per element of the input arrays, the runtime might generate so many threads that coordination of them will dominate the runtime. The overall strategy of “parallel, then sequential” is likely not the most efficient, either.

We can give a more refined strategy better exploiting the target architecture. For instance, modern CPUs support SIMD vector instructions which perform arithmetic operations in parallel on multiple scalar values in a single clock cycle. We would also like to trade off the amount of work that every thread performs with the number of threads launched to best exploit the available hardware cores while minimising the coordination overhead between the threads. The following expression groups the input data as vectors of 4 scalar elements before distributing the work across multiple threads using **mapPar**. Each thread then performs multiple sequential reductions in parallel on the vectorised data using SIMD multiplication and addition instructions, before a final sequential reduction adds up all remaining results.

$$\begin{aligned} &\text{reduceSeq } (+) \ 0 \ (\text{asScalar } (\text{join } (\text{mapPar} \\ &\quad (\text{mapSeq } (\lambda zs. \text{reduceSeq } (\lambda x \ a. \ (\text{fst } x * \text{snd } x) + a) \ 0 \ (\text{split } 2048 \ zs))) \\ &\quad (\text{split } (2048 * 64) \ (\text{zip } (\text{asVector } 4 \ xs) \ (\text{asVector } 4 \ ys)))))) \quad (3) \end{aligned}$$

Although this expression gives much more information about how to perform the computation in parallel, we have not left the functional paradigm. **mapPar**, **mapSeq**, and **reduceSeq** are just annotated variations of the standard **map** and **reduce** functions which still allows the straightforward mathematical reading of this expression. We can use equational reasoning to prove that this expression is semantically equivalent to (1). Equational reasoning can even be used to generate (2) and (3) from (1). Indeed Steuwer et al. [2016] have shown that stochastic search techniques are effective at automatically discovering parallelisation strategies that match hand-coded ones.

However, even with a specified parallelisation strategy we cannot execute this code directly. We need to translate the functional code to an imperative language like OpenMP or OpenCL in a way that preserves our chosen strategy. This paper presents a formal approach to solving this translation problem.

## 2.2 Strategy Preserving Translation to Imperative Code

What is the simplest way of converting a functional program to an imperative one? Starting with our naive “parallel map, sequential reduce” strategy of dot-product (2), we can turn it into an imperative program simply by assigning its result to an output variable *out*:

```
out := reduceSeq (+) 0 (mapPar (λx. fst x * snd x) (zip xs ys))
```

Unfortunately, this is not suitable for compilation targets like OpenMP or OpenCL. While assignment statements are the bread-and-butter of such languages, their expression languages certainly do not include such modern amenities as higher order *map* and *reduce* functions. To translate these away, we introduce a novel *acceptor-passing* translation  $\mathcal{A}(\llbracket E \rrbracket)_\delta(out)$ . The key idea is that for any expression *E* producing data of type  $\delta$ , the translation  $\mathcal{A}(\llbracket E \rrbracket)_\delta(out)$  is an imperative program that has the same effect as the assignment  $out := E$  and is free from higher-order combinators. This translation is mutually defined with a continuation passing translation  $C(\llbracket E \rrbracket)_\delta(C)$  that takes a parameterised command *C* that will consume the output, instead of taking an output variable.

The definition of the translation is given in Section 4.1. We introduce it here by example. Applied to our dot-product code, our translation first replaces the **reduceSeq** by a corresponding imperative combinator **reduceSeqI**. We will see below that **reduceSeqI** is straightforwardly implemented in terms of variable allocation and a for-loop.

$$\begin{aligned} & \mathcal{A}(\llbracket \text{reduceSeq } (+) 0 (\text{mapPar } (\lambda x. \text{fst } x * \text{snd } x) (\text{zip } xs \text{ ys})) \rrbracket)_{\text{num}}(out) \\ &= C(\llbracket \text{mapPar } (\lambda x. \text{fst } x * \text{snd } x) (\text{zip } xs \text{ ys}) \rrbracket)_{n.\text{num}}(\lambda x. \\ & \quad C(\llbracket 0 \rrbracket)_{\text{num}}(\lambda y. \text{reduceSeqI } n (\lambda x \ y \ o. \mathcal{A}(\llbracket x + y \rrbracket)_{\text{num}}(o)) \ y \ x (\lambda r. \mathcal{A}(\llbracket r \rrbracket)_{\text{num}}(out)))) \end{aligned}$$

The **mapPar** is translated, by the continuation passing translation, into allocation of a temporary array and an imperative **mapParI** combinator. As with **reduceSeqI**, the **mapParI** combinator is straightforwardly implementable in terms of a parallel for-loop. The operator **new**  $\delta$  *var* declares a new storage cell of type  $\delta$  named *var*, where storage cells are represented as pairs of an acceptor (i.e., “writer”, “l-value”) part *var.1* and an expression (i.e. “reader”, “r-value”) part *var.2*. Our language, which we introduce in Section 3, is a variant of Reynolds’ Idealised Algol [Reynolds 1978].

$$= \text{new } (n.\text{num}) (\lambda tmp. C(\llbracket \text{zip } xs \text{ ys} \rrbracket)_{n.(n \times \text{num})}(\lambda x. \text{mapParI } n (\lambda x \ o. \mathcal{A}(\llbracket \text{fst } x * \text{snd } x \rrbracket)_{\text{num}}(o)) \ x \ tmp.1); \\ C(\llbracket 0 \rrbracket)_{\text{num}}(\lambda y. \text{reduceSeqI } n (\lambda x \ y \ o. \mathcal{A}(\llbracket x + y \rrbracket)_{\text{num}}(o)) \ y \ tmp.2 (\lambda r. \mathcal{A}(\llbracket r \rrbracket)_{\text{num}}(out))))$$

Readers familiar with other translations of data parallel functional programs into imperative loops may be surprised at the allocation of a temporary array here. Typically, the compiler would be expected to automatically fuse the computation of the *map* into the translation of the *reduce*. However, this is precisely what we do *not* want from a predictable compilation process for parallelism. If fusion is desired, it is carried out before this translation is applied and directly encoded in the strategy of the functional program, as seen earlier in example (3). The parallelism strategy described by the functional code here precisely states “parallel map, followed by sequential reduce”.

Continuing the translation process, we substitute *out*, the arithmetic expressions and the **zip**, leaving two uses of the “intermediate-level” macros **mapParI** and **reduceSeqI**:

$$= \text{new } (n.\text{num}) (\lambda tmp. \text{mapParI } n (\lambda x \ o. o := \text{fst } x * \text{snd } x) (\text{zip } xs \text{ ys}) \ tmp.1; \\ \text{reduceSeqI } n (\lambda x \ y \ o. o := x + y) 0 \ tmp.2 (\lambda r. out := r))$$

These combinators are now replaced by parallel and sequential for-loops, which we describe in a further translation stage in Section 4.2.

$$= \text{new } (n.\text{num}) (\lambda tmp. \text{parfor } n \ tmp.1 (\lambda i \ o. o := \text{fst } (\text{idx } (\text{zip } xs \text{ ys}) \ i) * \text{snd } (\text{idx } (\text{zip } xs \text{ ys}) \ i)); \\ \text{new } \text{num } (\lambda accum. accum.1 := 0; \\ \text{for } n (\lambda i. accum.1 := accum.2 + \text{idx } tmp.2 \ i); \\ out := accum.2))$$

(4)

197 The sequential **for** loops of our intermediate language are standard; **for**  $n$  ( $\lambda i. b$ ) executes the body  
 198  $b$   $n$  times with iteration counter  $i$ . Parallel **parfor** loops are slightly more complex due to the way  
 199 they explicitly take a parameter (here named  $tmp.1$ ) that describes where to place the results of  
 200 each iteration in a data-race free way. We describe this fully in Section 3.3.

201 We are now left with an imperative program, albeit with a non-standard parallel-for construct  
 202 and complex data access expressions involving **fst**, **zip**, **idx** and so on. Our final translation to C and  
 203 OpenMP (Section 4.3) resolves these data layout expressions into explicit indexing computations:

```

204
205 1 float tmp[N];
206 2 #pragma omp parallel for
207 3 for (int i = 0; i < N; i += 1)
208 4   tmp[i] = xs[i] * ys[i];
209 5 float accum = 0.0;
210 6 for (int i = 0; i < N; i += 1)
211 7   accum = accum + tmp[i];
212 8 out = accum;

```

212 This resulting low level imperative code precisely implements the strategy “parallel map (lines  
 213 2–4), followed by sequential reduce (lines 5–8)” described by our functional expression (2).

214 Our naive dot-product code does not produce particularly complex code, but our translation  
 215 method scales to more detailed parallelism strategies. The alternative dot-product execution strategy  
 216 in (3), which rearranges the abstract *map* and *reduce* combinators in order to better exploit parallel  
 217 hardware, yields the following code using OpenMP pragmas for thread level parallelism (**omp**  
 218 **parallel for**) as well as for SIMD vector instructions (**omp simd**):

```

220 1 float tmp[N/(4*131072/64)];
221 2 #pragma omp parallel for
222 3 for (int i = 0; i < (N / (4 * 131072)); i += 1) {
223 4   for (int j = 0; j < 64; j += 1) {
224 5     float accum[4];
225 6     #pragma omp simd
226 7     for (int k = 0; k < 4; k += 1) { accum[k] = ((float[4]){0.0f})[k]; }
227 8     for (int l = 0; l < 2048; l += 1) {
228 9       #pragma omp simd
229 10      for (int k = 0; k < 4; k += 1) {
230 11        accum[k] = (xs[(524288 * i) + (8192 * j) + (4 * l) + k]
231 12          * ys[(524288 * i) + (8192 * j) + (4 * l) + k]) + accum[k];}
232 13      #pragma omp simd
233 14      for (int k = 0; k < 4; k += 1) { tmp[k + (4 * j) + (256 * i)] = accum[k];}
234 15    }
235 16  }
236 17 float accum = 0.0;
237 18 for (int i = 0; i < N/(4*131072/64); i += 1)
238 19   accum = accum + tmp[i];
239 20 out = accum;

```

237 As we shall see in Section 7, given a target-architecture optimised parallelisation strategy defined  
 238 in functional code, our translation process produces OpenMP code with performance on a par with  
 239 hand-written code.

240 Key to our translation methodology is a single intermediate language that can express pure  
 241 functional expressions and deterministic race free parallel imperative programs, and which is  
 242 amenable to formal reasoning. In the next section, we describe our language for this task, DPIA:  
 243 *Data Parallel Idealised Algol*.

244

245

### 3 DATA PARALLEL IDEALISED ALGOL

Our intermediate language for code generation is a dialect of Reynolds’s [1997] *Idealised Algol*, extended with interference control via a substructural type system [O’Hearn et al. 1999]. We saw some examples of DPIA in the previous section. We now highlight the major features of DPIA, and then give the formal presentation of its types (Section 3.1), typing (Section 3.2), and data parallel programming primitives (Section 3.3). We discuss the formal semantics of DPIA in Section 5.

Idealised Algol is built around Reynolds’ observation that we can use typed  $\lambda$ -calculus as a powerful macro language for imperative programming. The base types of the system are what would normally be the syntactic categories of a first order imperative language: `comm` is the type of commands, or statements; `exp` is the type of expressions, or r-values; and `acc` is the type of acceptors or l-values (entities that can be assigned to). The syntactic constructs of the language are expressed as constants in the language. For imperative programming, these include: `skip` : `comm` for doing nothing; `(;)` : `comm` × `comm` → `comm` for sequencing; and `(:=)` : `acc` × `exp` → `comm` for assignment. Declaration of new variables has a higher order type, `new` : (`acc` × `exp` → `comm`) → `comm`: we supply a command that requires an additional piece of storage to work with, and `new` allocates the variable, runs the command, and deallocates, following a block structured stack discipline.

The typing of `new` indicates a useful programming pattern in Idealised Algol. We can think of terms with types of the form  $Y \rightarrow (X \rightarrow \text{comm}) \rightarrow \text{comm}$  as implementing “objects” with interfaces  $X$ , using existing objects with interface  $Y$ . A canonical example is the counter object:

$$\begin{aligned} \text{counter} &: \text{acc} \rightarrow (\text{comm} \rightarrow \text{comm}) \rightarrow \text{comm} \\ \text{counter } o \text{ body} &= \mathbf{new} (\lambda \langle v_{wr}, v_{rd} \rangle. \text{body } (v_{wr} := v_{rd} + 1); o := v_{rd}) \end{aligned}$$

Given a client, `body`, and a place to write the final count to, `o`, `counter` allocates a new variable using `new` and hands `body` the restricted capability to only affect this variable by incrementing it. Once `body` has finished, the final value of the counter is written to `o`.

This implementation of `counter` uses  $\lambda$ -abstraction to hide its internal state from `body`, preventing any programming errors that would occur if `body` were able to interfere with `counter`’s internal state. However, in original Idealised Algol, `counter` can only encapsulate what it creates. It is possible for both `counter` and `body` to write to `o`, potentially producing confusing results. Reynolds’s [1978] *Syntactic Control of Interference* proposed a discipline that prevents such interference from taking place. Under interference control, `body` and `o` may only share identifiers that are used *passively* (i.e., read only). In use, interference control acts like Rust’s *borrow checker* [Klabnik and Nichols 2018]. In a term of the form `counter o (λinc. P); Q`, the variable `o` may not be used in `P`, except through `inc`. We may say that `inc`, the interface to the counter, has borrowed the resources owned by `o`. In the continuation command `Q`, this borrow ends and `o`’s resources may be used directly again.

We use this *safe resource reinterpretation* pattern multiple times in DPIA to provide abstracted access to low level parallelisation capabilities. At the lowest level, we have the `parfor` primitive:

$$\mathbf{parfor} : (n : \text{nat}) \rightarrow (t : \text{data}) \rightarrow \text{acc}[n.t] \rightarrow (\text{exp}[\text{idx}[n]] \rightarrow \text{acc}[t] \rightarrow_p \text{comm}) \rightarrow \text{comm}$$

The type of `parfor` is more elaborate than before to allow for datatype and size polymorphism, but the essence is this: the latter two of `parfor`’s arguments are an acceptor for writing results to, and a command for implementing the body of the loop. The same interference control discipline as above ensures that the body of the loop cannot write to the array chosen for output, except through the acceptor it is give, eliminating some potential data races. The `p` subscript on the type of the body indicates that this command must be passive, meaning that it cannot write to any state which it is not explicitly handed. Thus, it is possible to safely implement `parfor` as a parallel for loop. We use this same pattern to safely implement our intermediate level data parallel functions in Section 4.2 and our vectorisation and double buffering extensions in Section 6.

295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343

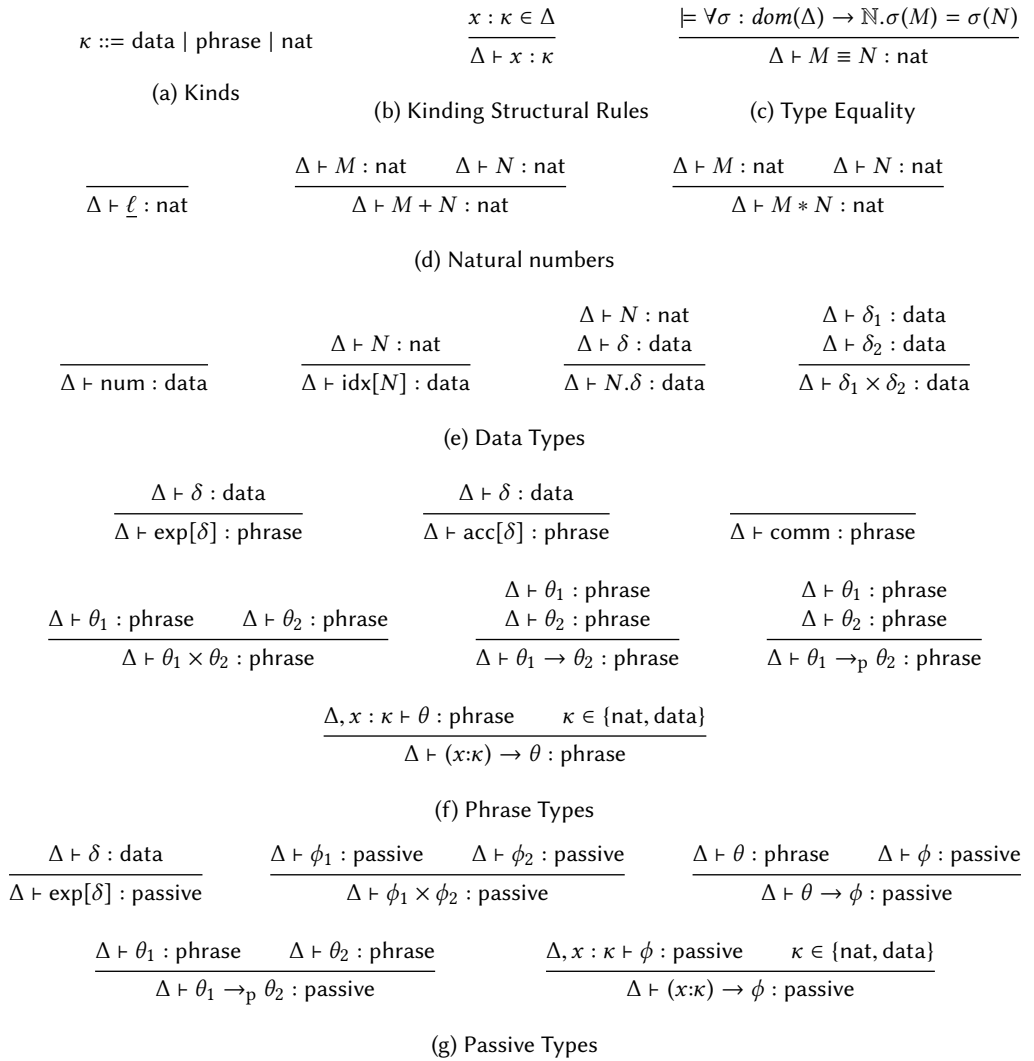


Fig. 1. Well-formed Types

### 3.1 The Types of DPIA

The type system of DPIA refines the *exp*, *acc*, *comm* *phrase types* described above by parameterising expressions and acceptors by *data types*, which classify data (integers, floats, arrays, *etc.*). To facilitate interference control, we identify a subset of phrase types which are *passive* (Section 3.1.2), i.e. essentially read-only, and so are safe to share across parallel threads. (We elaborate on what “essentially read-only” means in Section 3.1.2 and Section 3.2.)

**3.1.1 Kinding rules.** We extend SCIR with both data type and size polymorphism, so we need a kind system. Figure 1 presents the kinding rules for DPIA types. The kinds  $\kappa$  of DPIA include the major classifications into data types (*data*) and phrase types (*phrase*), along with the kind of type-level natural numbers (*nat*). Types may contain variables, so we use a kinding judgement

$\Delta \vdash \tau : \kappa$ , which states that type  $\tau$  has kind  $\kappa$  in kinding context  $\Delta$ . Figure 1b gives the variable rule that permits the use of type variables in well-kinded types. Figure 1d presents the rules for type-level natural numbers: either numeric literals  $\underline{\ell}$ , addition  $M + N$ , or multiplication  $M * N$  (where  $M$  and  $N$  range over terms of kind  $\text{nat}$ ). We write  $MN$  as syntactic sugar for  $M * N$ .

The rules for data types are presented in Figure 1e. We use  $\delta$  to range over data types. The base types are  $\text{num}$  for numbers; and a data type of array indexes  $\text{idx}[M]$ , parameterised by the size of the array. There are two compound types of data. For any data type  $\delta$  and natural number term  $M$ ,  $M.\delta$  is the data type of homogeneous arrays of  $\delta$ s of size  $M$ . We opt for a concise notation for array types as they are pervasive in data parallel programming. Heterogeneous compound data types (records) are built using the rule for  $\delta_1 \times \delta_2$ .

The phrase types of DPIA are given in Figure 1f. We use  $\theta$  to range over phrase types. For each data type  $\delta$ , there are phrase types  $\text{exp}[\delta]$  for *expression* phrases that produce data of type  $\delta$ , and  $\text{acc}[\delta]$  for *acceptor* phrases that consume data of type  $\delta$ . The  $\text{comm}$  phrase type classifies *command* phrases that may modify the store. Phrases that can be used in two different ways,  $\theta_1$  or  $\theta_2$ , are classified using the phrase product type  $\theta_1 \times \theta_2$ . This type is distinct from the *data* product type  $\delta_1 \times \delta_2$ : the data type represents a pair of data values; the phrase type represents an “interface” that offers two possible “methods”. The final three phrase types are all variants of parameterised phrase types. The phrase types  $\theta_1 \rightarrow \theta_2$  and  $\theta_1 \rightarrow_p \theta_2$  classify phrase functions. The  $p$  subscript denotes passive functions. The phrase type  $(x:\kappa) \rightarrow \theta$  classifies a phrase that is parameterised either by a data type or a natural number. We find it convenient to allow lowercase letters to range over both term variables and type variables.

The types of DPIA include arithmetic expressions, so we have a non trivial notion of equality between types, written  $\Delta \vdash \tau_1 \equiv \tau_2 : \kappa$ . The key type equality rule is given in Figure 1c: two arithmetic expressions are equal if they are equal as natural numbers for all interpretations ( $\sigma$ ) of their free variables. This equality is lifted to all other types by structural congruence.

**3.1.2 Passive Types.** Figure 1g identifies the subset of phrase types that classify passive phrases. The opposite of passive is *active*. We use  $\phi$  to range over passive phrase types. An expression phrase type  $\text{exp}[\delta]$  is always passive — phrases of this type can, by definition, only read the store. A compound phrase type is always passive if its component phrase types are all passive. Furthermore, a passive function type  $\theta_1 \rightarrow_p \theta_2$  is always passive, and a plain function type is passive whenever its return type is passive (irrespective of the argument type). Passive types are essentially read-only. The one exception whereby a phrase of passive type may modify the store is a passive function with active argument and return types. Such a function can only modify the part of the store addressable through the active phrase it is supplied with as an argument.

### 3.2 Typing Rules for DPIA

The typing judgement of DPIA follows the SCIR system of O’Hearn et al. [1999] in distinguishing between passive and active uses of identifiers. Our judgement also has a kinding context for size and data type polymorphism. The judgement form has the following structure:

$$\Delta \mid \Pi; \Gamma \vdash P : \theta$$

where  $\Delta$  is the kinding context,  $\Pi$  is a context of passively used identifiers,  $\Gamma$  is a context of actively used identifiers,  $P$  is a program phrase, and  $\theta$  is a phrase type. All the types in  $\Pi$  and  $\Gamma$  are phrase types well-kinded by  $\Delta$ . The phrase type  $\theta$  must also be well-kinded by  $\Delta$ . The order of entries does not matter. The contexts  $\Delta$  and  $\Pi$  allow contraction and weakening; context  $\Gamma$  does not.

The typing rules of DPIA are given in Figure 2. These rules define how variable phrases are formed, how parameterised and compound phrases are introduced and eliminated, and how passive



393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441

$$\begin{array}{c}
\frac{x : \theta \in \Gamma}{\Delta \mid \Pi; \Gamma \vdash x : \theta} \text{VAR} \qquad \frac{\Delta \mid \Pi; \Gamma \vdash P : \theta_1 \quad \Delta \vdash \theta_1 \equiv \theta_2 : \text{phrase}}{\Delta \mid \Pi; \Gamma \vdash P : \theta_2} \text{CONV} \qquad \frac{\text{prim} : \theta \in \text{PRIMITIVES}}{\Delta \mid \Pi; \Gamma \vdash \text{prim} : \theta} \text{PRIM} \\
\text{(a) Structural Rules} \\
\frac{\Delta \mid \Pi; \Gamma, x : \theta_1 \vdash P : \theta_2}{\Delta \mid \Pi; \Gamma \vdash \lambda x. P : \theta_1 \rightarrow \theta_2} \text{LAM} \qquad \frac{\Delta \mid \Pi; \Gamma_1 \vdash P : \theta_1 \rightarrow \theta_2 \quad \Delta \mid \Pi; \Gamma_2 \vdash Q : \theta_1}{\Delta \mid \Pi; \Gamma_1, \Gamma_2 \vdash P Q : \theta_2} \text{APP} \\
\frac{\Delta, x : \kappa \mid \Pi; \Gamma \vdash P : \theta \quad x \notin \text{ftv}(\Pi, \Gamma)}{\Delta \mid \Pi; \Gamma \vdash \Lambda x. P : (x:\kappa) \rightarrow \theta} \text{TLAM} \qquad \frac{\Delta \mid \Pi; \Gamma \vdash P : (x:\kappa) \rightarrow \theta \quad \Delta \vdash \tau : \kappa}{\Delta \mid \Pi; \Gamma \vdash P \tau : \theta[\tau/x]} \text{TAPP} \\
\frac{\Delta \mid \Pi; \Gamma \vdash P : \theta_1 \quad \Delta \mid \Pi; \Gamma \vdash Q : \theta_2}{\Delta \mid \Pi; \Gamma \vdash \langle P, Q \rangle : \theta_1 \times \theta_2} \text{PAIR} \qquad \frac{\Delta \mid \Pi; \Gamma \vdash P : \theta_1 \times \theta_2}{\Delta \mid \Pi; \Gamma \vdash P.i : \theta_i} \text{PROJ} \\
\text{(b) Introduction and Elimination Rules} \\
\frac{\Delta \mid \Pi, x : \theta; \Gamma \vdash P : \theta'}{\Delta \mid \Pi; \Gamma, x : \theta \vdash P : \theta'} \text{ACTIVATE} \qquad \frac{\Delta \mid \Pi; \Gamma, x : \theta \vdash P : \phi}{\Delta \mid \Pi, x : \theta; \Gamma \vdash P : \phi} \text{PASSIFY} \\
\frac{\Delta \mid \Pi; \cdot \vdash P : \theta_1 \rightarrow \theta_2}{\Delta \mid \Pi; \cdot \vdash P : \theta_1 \rightarrow_p \theta_2} \text{PROMOTE} \qquad \frac{\Delta \mid \Pi; \Gamma \vdash P : \theta_1 \rightarrow_p \theta_2}{\Delta \mid \Pi; \Gamma \vdash P : \theta_1 \rightarrow \theta_2} \text{DERELICT} \\
\text{(c) Active and Passive Phrase Rules}
\end{array}$$

Fig. 2. Typing Rules: Indexed Affine  $\lambda$ -Calculus with Passivity [O’Hearn et al. 1999]

and active types are managed. Any particular application of DPIA is specified by giving a collection of primitive phrases `PRIMITIVES`, each of which has a closed phrase type. We describe a collection for data parallel programming in Section 3.3, and extend it in Section 6.

Figure 2a presents the rule for forming variable phrases, implicit conversion between equal types, and the use of primitives. At point of use, all variables are considered to be used actively. If the final phrase type is passive, then an active use may be converted to a passive one by the `PASSIFY` rule. Primitives may be used in any context. Figure 2b presents the rules for parameterised and compound phrases. These are all standard typed  $\lambda$ -calculus style rules, except the use of separate contexts for a function and its arguments in the `APP` rule. This ensures that every function and its argument use non-interfering active resources, maintaining the invariant that distinct identifiers refer to non-interfering phrases. Note that we do not require separate contexts for the two parts of a compound phrase in the `PAIR` rule. Compound phrases offer two ways of interacting with the *same* underlying resource (as in the with (&) rule from Linear Logic).

Figure 2c describes how passive and active uses of variables are managed. The `ACTIVATE` rule allows any variable that has been used passively to be treated as if it were used actively. The `PASSIFY` rule allows active uses to be treated as passive, as long as the final phrase type is passive. The `PROMOTE` rule turns functions into passive functions, as long as they do not contain any free variables used actively. The `DERELICT` rule indicates that a passive function can always be seen as a normal function, if required.

DPIA’s *functional sub-language*. By inspection of the rules, we can see that if we restrict to phrase types constructed from `exp[\delta]`, functions, polymorphic functions, and tuples, then the constraints

on multiple uses of variables in DPIA cease to apply. Therefore, DPIA has a sub-language that has the same type system as a normal (non-substructural) typed  $\lambda$ -calculus with base types for numbers, arrays and tuples, and a limited form of polymorphism. When we introduce the functional primitives for DPIA in the next section, we will enrich this  $\lambda$ -calculus with arithmetic, array manipulators, and higher-order array combinators. It is this purely functional sub-language of DPIA that allows us to embed functional data parallel programs in a semantics preserving way.

### 3.3 Data Parallel Programming Primitives

The SCIR framework may be instantiated in arbitrary ways by specifying a collection of primitives. Here we instantiate it with typed primitive operations for data parallel programming, outlined in Figure 3. Our primitives fall into two principal categories: high-level functional primitives and low-level imperative primitives. Source programs include abstract high-level functional primitives **map** and **reduce**. Programs that are the input to our translation process are composed of concrete high-level functional primitives. All uses of **map** must be replaced by **mapSeq** (a sequential implementation of **map**) or **mapPar** (a parallel implementation of **map**), and all uses of **reduce** by **reduceSeq** (a sequential implementation of **reduce**). (An OpenCL backend might have other concrete implementations of **map** corresponding to the parallelisation hierarchy and a parallel implementation of **reduce** makes sense for some backends.) Uses of **mapPar**, **mapSeq**, and **reduceSeq** have no immediate counterpart in low-level languages for data-parallel computation. Our translation process converts these into low-level combinators (Section 4.1). A final lowering translation removes all functional primitives except arithmetic (Section 4.3).

As primitives are treated specially by the PRIM rule they can (with the aid of a little  $\eta$ -expansion) always be promoted to be passive. Thus, it is never necessary to annotate the arrows of a first-order primitive with a p subscript. The only such annotations that are necessary are those final arrows of function types occurring inside the type of a higher-order primitive that is required to be passive (in our case, only **parfor** and **mapParl**).

*Functional Primitives.* Figure 3a lists the type signatures of the primitives used for constructing purely functional data parallel programs. These fit into four groups. The first group consists of numeric literals ( $\ell$ ) and first-order operations on scalars (**negate**, (+), (-), (\*), (/)). The second group contains the abstract higher-order functional combinators for constructing array processing programs: **map** and **reduceSeq**. These have (the Idealised Algol renditions of) the standard types for map and reduce, extended with size information. The third group contains the concrete higher-order functional combinators for constructing array processing programs: **mapPar**, **mapSeq** and **reduceSeq**. The fourth group comprises functions for manipulating data layouts: **zip** joins two arrays of equal length into an array of pairs, **split** breaks a one dimensional array into a two dimensional array and **join** flattens a two dimensional array into a one dimensional array, **pair** constructs a pair, **fst** and **snd** deconstruct a pair, and **idx** looks up a value in an array. All of these primitives are data type indexed and those that operate on arrays are also size indexed.

*Example: dot-product.* The naive strategy for the dot-product example from Section 2 is written using the functional primitives like so, for input vectors  $xs$  and  $ys$  of length  $n$ , the only difference being that all of the size and data type information is described in detail (often we can infer these arguments and so in practice we often omit them):

```

reduceSeq n num num ( $\lambda x y. x + y$ ) 0 (mapPar n (num  $\times$  num) num
                                     ( $\lambda x. \mathbf{fst}$  num num  $x * \mathbf{snd}$  num num  $x$ )
                                     (zip n num num  $xs\ ys$ ))

```

491	$\underline{\ell}$	:	$\text{exp}[\text{num}]$
492	<b>negate</b>	:	$\text{exp}[\text{num}] \rightarrow \text{exp}[\text{num}]$
493	(+, *, /, -)	:	$\text{exp}[\text{num}] \times \text{exp}[\text{num}] \rightarrow \text{exp}[\text{num}]$
494	<b>map</b>	:	$(n : \text{nat}) \rightarrow (s \ t : \text{data}) \rightarrow (\text{exp}[s] \rightarrow \text{exp}[t]) \rightarrow \text{exp}[n.s] \rightarrow \text{exp}[n.t]$
495	<b>reduce</b>	:	$(n : \text{nat}) \rightarrow (s \ t : \text{data}) \rightarrow (\text{exp}[s] \rightarrow \text{exp}[t]) \rightarrow \text{exp}[t] \rightarrow \text{exp}[n.s] \rightarrow \text{exp}[t]$
496	<b>mapSeq</b>	:	$(n : \text{nat}) \rightarrow (s \ t : \text{data}) \rightarrow (\text{exp}[s] \rightarrow \text{exp}[t]) \rightarrow \text{exp}[n.s] \rightarrow \text{exp}[n.t]$
497	<b>mapPar</b>	:	$(n : \text{nat}) \rightarrow (s \ t : \text{data}) \rightarrow (\text{exp}[s] \rightarrow \text{exp}[t]) \rightarrow \text{exp}[n.s] \rightarrow \text{exp}[n.t]$
498	<b>reduceSeq</b>	:	$(n : \text{nat}) \rightarrow (s \ t : \text{data}) \rightarrow (\text{exp}[s] \rightarrow \text{exp}[t]) \rightarrow \text{exp}[t] \rightarrow \text{exp}[n.s] \rightarrow \text{exp}[t]$
499	<b>zip</b>	:	$(n : \text{nat}) \rightarrow (s \ t : \text{data}) \rightarrow \text{exp}[n.s] \rightarrow \text{exp}[n.t] \rightarrow \text{exp}[n.(s \times t)]$
500	<b>split</b>	:	$(n \ m : \text{nat}) \rightarrow (t : \text{data}) \rightarrow \text{exp}[nm.t] \rightarrow \text{exp}[m.n.t]$
501	<b>join</b>	:	$(n \ m : \text{nat}) \rightarrow (t : \text{data}) \rightarrow \text{exp}[n.m.t] \rightarrow \text{exp}[nm.t]$
502	<b>pair</b>	:	$(s \ t : \text{data}) \rightarrow \text{exp}[s] \rightarrow \text{exp}[t] \rightarrow \text{exp}[s \times t]$
503	<b>fst</b>	:	$(s \ t : \text{data}) \rightarrow \text{exp}[s \times t] \rightarrow \text{exp}[s]$
504	<b>snd</b>	:	$(s \ t : \text{data}) \rightarrow \text{exp}[s \times t] \rightarrow \text{exp}[t]$
505	<b>idx</b>	:	$(n : \text{nat}) \rightarrow (t : \text{data}) \rightarrow \text{exp}[n.t] \rightarrow \text{exp}[\text{idx}[n]] \rightarrow \text{exp}[t]$
506			(a) Functional primitives
507	(;)	:	$\text{comm} \times \text{comm} \rightarrow \text{comm}$
508	<b>skip</b>	:	$\text{comm}$
509	<b>new</b>	:	$(t : \text{data}) \rightarrow (\text{var}[t] \rightarrow \text{comm}) \rightarrow \text{comm}$ (where $\text{var}[t] = \text{acc}[t] \times \text{exp}[t] : \text{phrase}$ )
510	(:=)	:	$\text{acc}[\text{num}] \times \text{exp}[\text{num}] \rightarrow \text{comm}$
511	<b>for</b>	:	$(n : \text{nat}) \rightarrow (\text{exp}[\text{idx}[n]] \rightarrow \text{comm}) \rightarrow \text{comm}$
512	<b>parfor</b>	:	$(n : \text{nat}) \rightarrow (t : \text{data}) \rightarrow \text{acc}[n.t] \rightarrow (\text{exp}[\text{idx}[n]] \rightarrow \text{acc}[t] \rightarrow_p \text{comm}) \rightarrow \text{comm}$
513	<b>zipAcc<sub>1</sub></b>	:	$(n : \text{nat}) \rightarrow (s \ t : \text{data}) \rightarrow \text{acc}[n.s \times t] \rightarrow \text{acc}[n.s]$
514	<b>zipAcc<sub>2</sub></b>	:	$(n : \text{nat}) \rightarrow (s \ t : \text{data}) \rightarrow \text{acc}[n.s \times t] \rightarrow \text{acc}[n.t]$
515	<b>splitAcc</b>	:	$(n \ m : \text{nat}) \rightarrow (t : \text{data}) \rightarrow \text{acc}[m.n.t] \rightarrow \text{acc}[nm.t]$
516	<b>joinAcc</b>	:	$(n \ m : \text{nat}) \rightarrow (t : \text{data}) \rightarrow \text{acc}[nm.t] \rightarrow \text{acc}[n.m.t]$
517	<b>pairAcc<sub>1</sub></b>	:	$(s \ t : \text{data}) \rightarrow \text{acc}[s \times t] \rightarrow \text{acc}[s]$
518	<b>pairAcc<sub>2</sub></b>	:	$(s \ t : \text{data}) \rightarrow \text{acc}[s \times t] \rightarrow \text{acc}[t]$
519	<b>idxAcc</b>	:	$(n : \text{nat}) \rightarrow (t : \text{data}) \rightarrow \text{acc}[n.t] \rightarrow \text{exp}[\text{idx}[n]] \rightarrow \text{acc}[t]$
520			(b) Imperative primitives
521			
522	<b>mapSeqI</b>	:	$(n : \text{nat}) \rightarrow (s \ t : \text{data}) \rightarrow (\text{exp}[s] \rightarrow \text{acc}[t] \rightarrow \text{comm}) \rightarrow \text{exp}[n.s] \rightarrow \text{acc}[n.t] \rightarrow \text{comm}$
523	<b>mapParI</b>	:	$(n : \text{nat}) \rightarrow (s \ t : \text{data}) \rightarrow (\text{exp}[s] \rightarrow \text{acc}[t] \rightarrow_p \text{comm}) \rightarrow \text{exp}[n.s] \rightarrow \text{acc}[n.t] \rightarrow \text{comm}$
524	<b>reduceSeqI</b>	:	$(n : \text{nat}) \rightarrow (s \ t : \text{data}) \rightarrow (\text{exp}[s] \rightarrow \text{exp}[t] \rightarrow \text{acc}[t] \rightarrow \text{comm}) \rightarrow$ $\text{exp}[t] \rightarrow \text{exp}[n.s] \rightarrow (\text{exp}[t] \rightarrow \text{comm}) \rightarrow \text{comm}$
525			(c) Intermediate imperative macros
526			
527			
528			
529			
530			
531			
532			
533			
534			
535			
536			
537			
538			
539			

Fig. 3. Data Parallel Programming Primitives, Functional and Imperative

Likewise, the more sophisticated strategy of dot-product from Section 2 with nested **splits** and **joins** can be expressed with detailed size and type information throughout.

*Imperative Primitives.* Figure 3b gives the type signatures for the imperative primitives. These are split into two groups. The first group includes the standard Idealised Algol combinators that turn DPIA into an imperative programming language: (;) sequences commands, **skip** is the command that does nothing; **new**  $\delta$  allocates a new mutable variable on the stack, where a variable is a pair of an acceptor and an expression; and (:=) assigns the value of an expression to an acceptor.

For-loops are constructed by the combinators **for** and **parfor**. Sequential for-loops **for**  $N B$  take a number of iterations  $N$  and a loop body  $B$ , a command parameterised by the iteration number. Parallel for-loops **parfor**  $N \delta A B$  differ from normal for-loops by taking an additional acceptor argument  $A : \text{acc}[B.\delta]$  that is used for the output of each iteration. As described in the introduction to this section, this allows an implementation of **parfor** that can safely be implemented using actual parallelism. We will see how the acceptor-transforming behaviour of our **parfor** primitive is translated into a normal, potentially racy, parallel for loop in OpenMP C in Section 4.3. Formally, newly allocated variables are zero initialised (and pointwise zero initialised for compound data), but in our implementation we typically optimise away the initialisation. In particular, it is never necessary to initialise dynamic memory allocations that are introduced by the translation of the functional primitives into imperative code as all dynamically allocated memory is always written to before being read.

The second group of imperative primitives include the acceptor variants of the **zip**, **split**, **join**, **pair** and **idx** functional primitives. The acceptor primitives transform acceptors of compound data into acceptors of their components. They are used to funnel data into the correct positions in the imperative translations of functional programs. In the final translation to parallel C code, described in Section 4.3, all acceptor phrases are translated into l-values with explicit index computations.

*Intermediate Imperative Macros.* Figure 3c gives type signatures for the intermediate imperative counterparts of **mapPar**, **mapSeq** and **reduceSeq**. These macros will be used in our translation from higher-order functional programs to higher-order imperative programs in Section 4.1. In the second stage of the translation they will be substituted by implementations in terms of variable allocation and for-loops (Section 4.2).

## 4 FROM FUNCTIONAL TO IMPERATIVE

In this section we spell out the three stages of translation from higher-order functional array programs to parallel C-like code (as sketched in Section 2). First, the higher-order functional combinators **mapPar**, **mapSeq** and **reduceSeq** are translated into the higher-order imperative macros **mapParI**, **mapSeqI** and **reduceSeqI** using an acceptor-passing translation defined mutually recursively with a continuation-passing translation (Section 4.1). Second, the higher-order imperative macros are expanded into for-loops (Section 4.2). Finally, low-level parallel pseudo-C code is produced by translating expression and acceptor phrases into indexing operations (Section 4.3). We prove the correctness of the first two stages of our translation in Section 5.

### 4.1 Translation Stage I: Higher-order Functional to Higher-order Imperative

The goal of the first stage of the translation is to take a phrase  $E : \text{exp}[\delta]$ , constructed from the functional primitives in Figure 3a, and an acceptor  $out : \text{acc}[\delta]$  and to produce a comm phrase that has the same semantics as the command

$$out :=_{\delta} E$$

where  $(:=_{\delta})$  is an assignment operator for non-base types defined by induction on  $\delta$  below. The resulting program will be an imperative program that acts as if we could compute the functional expression in one go and assign it to the output acceptor. Since our compilation targets know nothing of higher-order functional combinators like **mapPar**, **mapSeq** and **reduceSeq** they will have to be translated away. We do not use any of the traditional methods for compiling higher-order functions, such as closure conversion [Steele 1978] or defunctionalisation [Reynolds 1998]. Instead, we rely on the whole-program nature of our translation, our lack of recursion, and the special form of our functional primitives. Specifically, we are relying on a version of Gentzen's subformula

589	$\mathcal{A}(x)_\delta(A)$	=	$A :=_\delta x$
590	$\mathcal{A}(\ell)_{\text{num}}(A)$	=	$A :=_\ell$
591	$\mathcal{A}(\mathbf{negate} E)_{\text{num}}(A)$	=	$C(E)_{\text{num}}(\lambda x. A := \mathbf{negate} x)$
592	$\mathcal{A}(E_1 + E_2)_{\text{num}}(A)$	=	$C(E_1)_{\text{num}}(\lambda x. C(E_2)_{\text{num}}(\lambda y. A := x + y))$
593	$\mathcal{A}(\mathbf{mapPar} N \delta_1 \delta_2 F E)_{N.\delta_2}(A)$	=	$C(E)_{N.\delta_1}(\lambda x. \mathbf{mapPar} N \delta_1 \delta_2 (\lambda x o. \mathcal{A}(F x)_{\delta_2}(o)) x A)$
594	$\mathcal{A}(\mathbf{mapSeq} N \delta_1 \delta_2 F E)_{N.\delta_2}(A)$	=	$C(E)_{N.\delta_1}(\lambda x. \mathbf{mapSeq} N \delta_1 \delta_2 (\lambda x o. \mathcal{A}(F x)_{\delta_2}(o)) x A)$
595	$\mathcal{A}(\mathbf{reduceSeq} N \delta_1 \delta_2 F I E)_{\delta_2}(A)$	=	$C(E)_{N.\delta_1}(\lambda x. C(I)_{\delta_2}(\lambda y. \mathbf{reduceSeq} N \delta_1 \delta_2 (\lambda x y o. \mathcal{A}(F x y)_{\delta_2}(o)) y x (\lambda r. \mathcal{A}(r)(A))))$
596			
597	$\mathcal{A}(\mathbf{zip} N \delta_1 \delta_2 E_1 E_2)_{N.\delta_1 \times \delta_2}(A)$	=	$\mathcal{A}(E_1)_{N.\delta_1}(\mathbf{zipAcc}_1 N \delta_1 \delta_2 A); \mathcal{A}(E_2)_{N.\delta_2}(\mathbf{zipAcc}_2 N \delta_1 \delta_2 A)$
598	$\mathcal{A}(\mathbf{split} N M \delta E)_{N.M.\delta}(A)$	=	$\mathcal{A}(E)_{N.M.\delta}(\mathbf{splitAcc} N M \delta A)$
599	$\mathcal{A}(\mathbf{join} N M \delta E)_{N.M.\delta}(A)$	=	$\mathcal{A}(E)_{N.M.\delta}(\mathbf{joinAcc} N M \delta A)$
600	$\mathcal{A}(\mathbf{pair} \delta_1 \delta_2 E_1 E_2)_{\delta_1 \times \delta_2}(A)$	=	$\mathcal{A}(E_1)_{\delta_1}(\mathbf{pairAcc}_1 \delta_1 \delta_2 A); \mathcal{A}(E_2)_{\delta_2}(\mathbf{pairAcc}_2 \delta_1 \delta_2 A)$
601	$\mathcal{A}(\mathbf{fst} \delta_1 \delta_2 E)_{\delta_1}(A)$	=	$C(E)_{\delta_1 \times \delta_2}(\lambda x. A :=_{\delta_1} \mathbf{fst} \delta_1 \delta_2 x)$
602	$\mathcal{A}(\mathbf{snd} \delta_1 \delta_2 E)_{\delta_2}(A)$	=	$C(E)_{\delta_1 \times \delta_2}(\lambda x. A :=_{\delta_2} \mathbf{snd} \delta_1 \delta_2 x)$
603			
604		(a) Acceptor-passing Translation	
605	$C(x)_\delta(C)$	=	$C(x)$
606	$C(\ell)_{\text{num}}(C)$	=	$C(\ell)$
607	$C(\mathbf{negate} E)_{\text{num}}(C)$	=	$C(E)_{\text{num}}(\lambda x. C(\mathbf{negate} x))$
608	$C(E_1 + E_2)_{\text{num}}(C)$	=	$C(E_1)_{\text{num}}(\lambda x. C(E_2)_{\text{num}}(\lambda y. C(x + y)))$
609	$C(\mathbf{mapPar} N \delta_1 \delta_2 F E)_{N.\delta_2}(C)$	=	$\mathbf{new} (N.\delta_2) (\lambda tmp. \mathcal{A}(\mathbf{mapPar} N \delta_1 \delta_2 F E)_{N.\delta_2}(tmp.2); C(tmp.1))$
610	$C(\mathbf{mapSeq} N \delta_1 \delta_2 F E)_{N.\delta_2}(C)$	=	$\mathbf{new} (N.\delta_2) (\lambda tmp. \mathcal{A}(\mathbf{mapSeq} N \delta_1 \delta_2 F E)_{N.\delta_2}(tmp.2); C(tmp.1))$
611	$C(\mathbf{reduceSeq} N \delta_1 \delta_2 F I E)_{\delta_2}(C)$	=	$C(E)_{N.\delta_1}(\lambda x. C(I)_{\delta_2}(\lambda y. \mathbf{reduceSeq} N \delta_1 \delta_2 (\lambda x y o. \mathcal{A}(F x y)_{\delta_2}(o)) y x C))$
612	$C(\mathbf{zip} N \delta_1 \delta_2 E_1 E_2)_{N.\delta_1 \times \delta_2}(C)$	=	$C(E_1)_{N.\delta_1}(\lambda x. C(E_2)_{N.\delta_2}(\lambda y. C(\mathbf{zip} N \delta_1 \delta_2 x y)))$
613	$C(\mathbf{split} N M \delta E)_{N.M.\delta}(C)$	=	$C(E)_{N.M.\delta}(\lambda x. C(\mathbf{split} N M \delta x))$
614	$C(\mathbf{join} N M \delta E)_{N.M.\delta}(C)$	=	$C(E)_{N.M.\delta}(\lambda x. C(\mathbf{join} N M \delta x))$
615	$C(\mathbf{pair} \delta_1 \delta_2 E_1 E_2)_{\delta_1 \times \delta_2}(C)$	=	$C(E_1)_{\delta_1}(\lambda x. C(E_2)_{\delta_2}(\lambda y. C(\mathbf{pair} \delta_1 \delta_2 x y)))$
616	$C(\mathbf{fst} \delta_1 \delta_2 E)_{\delta_1}(C)$	=	$C(E)_{\delta_1 \times \delta_2}(\lambda x. C(\mathbf{fst} \delta_1 \delta_2 x))$
617	$C(\mathbf{snd} \delta_1 \delta_2 E)_{\delta_2}(C)$	=	$C(E)_{\delta_1 \times \delta_2}(\lambda x. C(\mathbf{snd} \delta_1 \delta_2 x))$
618			
619		(b) Continuation-passing Translation	
620			

Fig. 4. Acceptor and Continuation-passing Translations

principle (identified by Gentzen [1935] and named by Prawitz [1965]). Our approach is reminiscent of that of Najd et al. [2016] who use quotation and normalisation, making essential use of the subformula principle, to embed domain-specific languages. An obvious difference with our work is that rather than stratifying a language into a host functional language and a quoted functional language, we seamlessly combine a functional and an imperative language.

We have already mentioned the use of assignment at compound data types. This is defined by:

$$\begin{aligned}
 A :=_{\text{num}} E &= A := E \\
 A :=_{N.\delta} E &= \mathbf{mapPar} N \delta \delta (\lambda x a.a :=_\delta x) E A \\
 A :=_{\delta_1 \times \delta_2} E &= \mathbf{pairAcc}_1 A :=_{\delta_1} \mathbf{fst} E; \mathbf{pairAcc}_2 A :=_{\delta_2} \mathbf{snd} E
 \end{aligned}$$

The translation of functional expressions to imperative code is accomplished by two type-directed mutually defined translations: the acceptor-passing translation  $\mathcal{A}(-)_\delta$  (Figure 4a) and the continuation-passing translation  $C(-)_\delta$  (Figure 4b). The acceptor-passing translation takes a data type  $\delta$ , an expression of type  $\text{exp}[\delta]$  and an acceptor of type  $\text{acc}[\delta]$ , and produces a comm phrase.

Likewise, the continuation-passing translation takes a data type  $\delta$ , an expression of type  $\text{exp}[\delta]$  and a continuation of type  $\text{exp}[\delta] \rightarrow \text{comm}$ , and produces a  $\text{comm}$  phrase.

It is straightforward to see by inspection, and using the fact that weakening is admissible in DPIA, that the two translations are type-preserving:

THEOREM 4.1.

- (1) If  $\Delta \mid \Pi; \Gamma_1 \vdash E : \text{exp}[\delta]$  and  $\Delta \mid \Pi; \Gamma_2 \vdash A : \text{acc}[\delta]$  then  $\Delta \mid \Pi; \Gamma_1, \Gamma_2 \vdash \mathcal{A}(\langle E \rangle_\delta(A)) : \text{comm}$ .
- (2) If  $\Delta \mid \Pi; \Gamma_1 \vdash E : \text{exp}[\delta]$  and  $\Delta \mid \Pi; \Gamma_2 \vdash C : \text{exp}[\delta] \rightarrow \text{comm}$  then  $\Delta \mid \Pi; \Gamma_1, \Gamma_2 \vdash C(\langle E \rangle_\delta(C)) : \text{comm}$ .

It is also important that these translations satisfy the following equivalences.

$$\mathcal{A}(\langle E \rangle_\delta(A)) \simeq A :=_\delta E \qquad C(\langle E \rangle_\delta(C)) \simeq C(E)$$

We define observational equivalence ( $\simeq$ ) for DPIA and establish these particular equivalences in Section 5. Ultimately, our goal is to compute the result of  $\mathcal{A}(\langle E \rangle_\delta(\text{out}))$ . It might appear that we could dispense with the acceptor-passing translation and simply use  $C(\langle E \rangle_\delta(\lambda x. \text{out} :=_\delta x))$ . However, this would create unnecessary temporary storage, violating our desire for an efficient translation. There are clear similarities between our mutually-defined translations and tail-recursive one-pass CPS translations that avoid producing unnecessary administrative redexes [Danvy et al. 2007].

The clauses for both translations split into four groups. The first group consists only of the clause for translating functional expression phrases that are identifiers  $x$ . In the acceptor-passing case, we defer to the generalised assignment defined above; in the continuation-passing case, we simply apply the continuation to the variable. The second group handles first-order operations on numeric data. In all cases, we defer to the continuation-passing translation for sub-expressions, with appropriate continuations. The third group is the most interesting: the translations of the higher order **mapPar**, **mapSeq** and **reduceSeq** primitives. For **mapPar/mapSeq**: in the acceptor-passing case, we can immediately translate to **mapParl/mapSeq** which already takes an acceptor to place its output into; in the continuation-passing case, we must create a temporary array as storage, invoke **mapParl/mapSeqI**, and then let the continuation read from the temporary array. This indirection is required because we do not know what random access strategy the continuation  $C$  will use to read the array it is given. For **reduceSeq**, in both cases we translate to the **reduceSeqI** combinator. The fourth group handles the translations of the functional data layout combinators. The CPS translation, passes them straight through; they are handled by the final translation to low-level C-like code in Section 4.3. The acceptor-passing translation translates the combinators that construct data into corresponding acceptors. In the **fst** and **snd** cases, which project out data, we defer to the continuation-passing translation. In practice, the case of a projection in tail position rarely arises, since it corresponds to disposal of part of the overall computation.

*Limits of Strategy-Preservation.* Our translation is almost, but not quite, strategy-preserving. There are two places where the strategy is hard-wired. First, array assignments are always computed in parallel using **mapParl**. Second, the CPS translations of **mapPar** and **mapSeq** each introduce a temporary array, and there is no choice as to where in the memory hierarchy the array is allocated. If we were targeting a heterogeneous backend such as a GPU then we would likely require more fine-grained control over these choices. We defer study of such systems as well as a formal treatment of precisely what strategy-preservation means to future work.

## 4.2 Translation Stage II: Higher-order Imperative to For-loops

The next stage in the translation expands the intermediate level imperative macros **mapParl**, **mapSeqI** and **reduceSeqI** as lower-level implementations in terms of (parallel and sequential) for-loops. This is accomplished by substitution and  $\beta$ -reduction (DPIA includes full  $\beta\eta$  reasoning

principles). The combinator **mapParl** is implemented as a parallel loop:

$$\mathbf{mapParl} = \Lambda n s t. \lambda f xs out. \mathbf{parfor} n t out (\lambda i o. f (\mathbf{idx} n s xs i) o)$$

**mapSeqI** is implemented as a sequential loop:

$$\mathbf{mapSeqI} = \Lambda n s t. \lambda f xs out. \mathbf{for} n (\lambda i. f (\mathbf{idx} n s xs i) (\mathbf{idxAcc} n t out i))$$

The implementation of **reduceI** is more complex, involving the allocation of a temporary variable to store the accumulated value during the reduction. In this case, the for loop is sequential, since the semantics of reduction demands that we visit each element in turn:

$$\begin{aligned} \mathbf{reduceSeqI} = \Lambda n s t. \lambda f init xs c. \mathbf{new} t (\lambda acc. acc.2 :=_s init; \\ \mathbf{for} n (\lambda i. f (\mathbf{idx} n s xs i) (acc.2) (acc.1)); \\ c (acc.1)) \end{aligned}$$

### 4.3 Translation Stage III: For-loops to Parallel Pseudo-C

After performing the translation steps in the previous two sections, we have generated a command phrase that does not use the higher-order functional combinators **mapPar**, **mapSeq** and **reduceSeq**, but still contains uses of the data layout combinators **zip**, **split** etc., and their acceptor counterparts. We now define a translation to OpenMP C, an extension of the C language adding parallel for loops. This translation resolves the data layout expressions into explicit indexing.

The translation in Figure 5 is split into three parts: the translation of DPIA commands into OpenMP C statements, the translation of acceptors into l-values, and the translation of expressions into r-values. (Recall the analogy made in Section 3.1 between the phrase types of DPIA and syntactic categories in an imperative language; we are now reaping the rewards of Reynolds' careful design of Idealised Algol.) We assume that the input to the translation process is in  $\beta\eta$ -normal form, so all **new**-block and loop bodies are fully expanded.

The translation of commands in Figure 5a straightforwardly translates each DPIA command into the corresponding C statement. The translation is parameterised by an environment  $\eta$  that maps from DPIA identifiers to C variable names. There is a small discrepancy that we overcome in that semicolons are statement terminators in C, but separators in DPIA, and that doing nothing useful is unremarkable in a C program, but is explicitly written **skip** in DPIA. The translation of assignment relies on the translations for acceptors and expression, which we define below. Variable allocation is translated using a  $\{ \dots \}$  block to limit the scope. We omit initialisation of the new variable because we know by inspection of our previous translation steps that newly allocated variables will always be completely initialised before reading. Note how we explicitly substitute a pair of identifiers for the acceptor and expression parts of the variable in DPIA, but that these both refer to the same C variable in the extended environment. Translation of variable allocation makes use of generator  $\text{CODEGEN}_{\text{data}}(\delta, v)$  that generates the appropriate C-style variable declaration for the data type  $\delta$ . Since C does not have anonymous tuple types, this entails on-the-fly generation of the appropriate **struct** definitions. DPIA **for** loops are translated into C for loops, DPIA **parfor** loops are translated into pseudo-parallel-for loops. In the body of the **parfor** loop, we substitute in an **idxAcc** phrase which will be resolved later by the translation of acceptors. The variable names introduced by translating **new**, **for** and **parfor** are all assumed to be fresh.

The translation of acceptors (Figure 5b) is parameterised by an environment  $\eta$ , as for commands, as well as a path  $ps$ , consisting of a list of C-expressions of type `int` denoting array indexes, and **struct** fields,  $.x1$ ,  $.x2$ , denoting projections from pairs. The path must always agree with the type of the acceptor being translated. During command translation, all acceptors being translated have type `num` so the access path starts empty. The acceptor translation clauses in Figure 5b all proceed

```

736 CODEGENcomm(skip,  $\eta$ ) = /* skip */
737 CODEGENcomm( $P_1; P_2, \eta$ ) = CODEGENcomm( $P_1, \eta$ ) CODEGENcomm( $P_2, \eta$ )
738 CODEGENcomm( $A := E, \eta$ ) = CODEGENacc[num]( $A, \eta, []$ ) = CODEGENexp[num]( $E, \eta, []$ );
739 CODEGENcomm(new  $\delta$  ( $\lambda v. P$ ),  $\eta$ ) = { CODEGENdata( $\delta$ ) v;
740 CODEGENcomm( $P[(v_a, v_e)/v], \eta[v_a \mapsto v, v_e \mapsto v]$ ) }
741 CODEGENcomm(for  $n$  ( $\lambda i. P$ ),  $\eta$ ) = for(int i = 0; i < n; i += 1) {
742 CODEGENcomm( $P, \eta[i \mapsto i]$ ) }
743 CODEGENcomm(parfor  $n$   $\delta$   $A$  ( $\lambda i o. P$ )  $E, \eta$ ) = #pragma omp parallel for
744 for(int i = 0; i < n; i += 1) {
745 CODEGENcomm( $P[\mathbf{idxAcc} n \delta A i/o], \eta[i \mapsto i]$ ) }

```

## (a) DPIA commands to OpenMP C statements

```

748 CODEGENacc[ $\delta$ ]( $x, \eta, ps$ ) =  $\eta(x)$ (reverse  $ps$ )
749 CODEGENacc[ $\delta$ ](idxAcc  $N \delta A I, \eta, ps$ ) = CODEGENacc[N. $\delta$ ]( $A, \eta,$ 
750 CODEGENexp[idx[N]]( $I, \eta, []$ ) ::  $ps$ ))
751 CODEGENacc[N.M. $\delta$ ](splitAcc  $N M \delta A, \eta, i :: ps$ ) = CODEGENacc[N.M. $\delta$ ]( $A, \eta, i/M :: i \% M :: ps$ )
752 CODEGENacc[N.M. $\delta$ ](joinAcc  $n M \delta A, \eta, i :: j :: ps$ ) = CODEGENacc[N.M. $\delta$ ]( $A, \eta, i * M + j :: ps$ )
753 CODEGENacc[ $\delta_1$ ](pairAcc1  $\delta_1 \delta_2 A, \eta, ps$ ) = CODEGENacc[ $\delta_1 \times \delta_2$ ]( $A, \eta, .x1 :: ps$ )
754 CODEGENacc[ $\delta_2$ ](pairAcc2  $\delta_1 \delta_2 A, \eta, ps$ ) = CODEGENacc[ $\delta_1 \times \delta_2$ ]( $A, \eta, .x2 :: ps$ )
755 CODEGENacc[N. $\delta_1$ ](zipAcc1  $N \delta_1 \delta_2 A, \eta, i :: ps$ ) = CODEGENacc[N. $(\delta_1 \times \delta_2)$ ]( $A, \eta, i :: .x1 :: ps$ )
756 CODEGENacc[N. $\delta_2$ ](zipAcc2  $N \delta_1 \delta_2 A, \eta, i :: ps$ ) = CODEGENacc[N. $(\delta_1 \times \delta_2)$ ]( $A, \eta, i :: .x2 :: ps$ )

```

## (b) DPIA acceptors to C l-values

```

759 CODEGENexp[ $\delta$ ]( $x, \eta, ps$ ) =  $\eta(x)$ (reverse  $ps$ )
760 CODEGENexp[num]( $\ell, \eta, []$ ) =  $\ell$ 
761 CODEGENexp[num](negate  $E, \eta, []$ ) = (- CODEGENexp[num]( $E, \eta, []$ ))
762 CODEGENexp[num]( $E_1 + E_2, \eta, []$ ) = (CODEGENexp[num]( $E_1, \eta, []$ )
763 + CODEGENexp[num]( $E_2, \eta, []$ ))
764 CODEGENexp[N. $(\delta_1 \times \delta_2)$ ](zip  $N \delta_1 \delta_2 E_1 E_2, \eta, i :: .xj :: ps$ ) = CODEGENexp[N. $\delta_j$ ]( $E_j, \eta, i :: ps$ )
765 CODEGENexp[M.N. $\delta$ ](split  $N M \delta E, \eta, i :: j :: ps$ ) = CODEGENexp[M.N. $\delta$ ]( $E, \eta, i * n + j :: ps$ )
766 CODEGENexp[M.N. $\delta$ ](join  $N M \delta E, \eta, i :: ps$ ) = CODEGENexp[M.N. $\delta$ ]( $E, \eta, i / n :: i \% n :: ps$ )
767 CODEGENexp[ $\delta_1 \times \delta_2$ ](pair  $\delta_1 \delta_2 E_1 E_2, \eta, .xj :: ps$ ) = CODEGENexp[ $\delta_j$ ]( $E_j, \eta, ps$ )
768 CODEGENexp[ $\delta_1$ ](fst  $\delta_1 \delta_2 E, \eta, ps$ ) = CODEGENexp[ $\delta_1 \times \delta_2$ ]( $E, \eta, .x1 :: ps$ )
769 CODEGENexp[ $\delta_2$ ](snd  $\delta_1 \delta_2 E, \eta, ps$ ) = CODEGENexp[ $\delta_1 \times \delta_2$ ]( $E, \eta, .x2 :: ps$ )
770 CODEGENexp[ $\delta$ ](idx  $N \delta E I, \eta, ps$ ) = CODEGENexp[ $N. $\delta$ ]( $E, \eta,$ 
771 CODEGENexp[idx[N]]( $I, \eta, []$ ) ::  $ps$ )$ 
```

## (c) DPIA expressions to C r-values

Fig. 5. Translation of Purely Imperative DPIA to parallel OpenMP C

by manipulating the access path appropriately until an identifier is reached. At this point, the DPIA identifier is replaced with its corresponding C variable and the access path is appended.

Translation of expressions (Figure 5c) is parameterised similarly to the acceptor translation, and contains similar clauses for all the data layout combinators. Expressions also include literals and arithmetic expressions, which are translated to the corresponding notion in C.



*Example.* We demonstrate how the translation to C works by applying it to the **parfor** loop in the translation (4) of the simple dot product in Section 2. We use the environment  $\eta = [out \mapsto out, xs \mapsto xs, ys \mapsto ys]$ . The command translation translates, in two steps, the **parfor** loop and the assignment, substituting in the indexing acceptor for the acceptor in the loop body:

```

785 CODEGENcomm(parfor  $\ell$  out ( $\lambda i o.o := \text{fst}(\text{idx}(\text{zip } xs \text{ } ys) i) * \text{snd}(\text{idx}(\text{zip } xs \text{ } ys) i)$ ),  $\eta$ )
786 = #pragma omp parallel for
787 for(int i = 0; i <  $\ell$ ; i+=1) {
788   CODEGENcomm(idxAcc out i :=  $\text{fst}(\text{idx}(\text{zip } xs \text{ } ys) i) * \text{snd}(\text{idx}(\text{zip } xs \text{ } ys) i)$ ),  $\eta[i \mapsto i]$ )
789 }
790 = #pragma omp parallel for
791 for(int i = 0; i <  $\ell$ ; i+=1) {
792   CODEGENacc[num](idxAcc out i,  $\eta[i \mapsto i]$ , []) =
793   CODEGENexp[num]( $\text{fst}(\text{idx}(\text{zip } xs \text{ } ys) i) * \text{snd}(\text{idx}(\text{zip } xs \text{ } ys) i)$ ),  $\eta[i \mapsto i]$ , []);
794 }
795 }

```

The acceptor part of the assignment is translated as follows:

```

800 CODEGENacc[num](idxAcc out i,  $\eta[i \mapsto i]$ , []) = CODEGENacc[ $\ell$ .num](out,  $\eta[i \mapsto i]$ , [i])
801 = out[i]

```

The expression part of the assignment is translated as follows, where we only spell out the left-hand side of the multiplication in detail; the right hand side is similar.

```

802 CODEGENexp[num]( $\text{fst}(\text{idx}(\text{zip } xs \text{ } ys) i)$ ,  $\eta[i \mapsto i]$ , [])
803 = CODEGENexp[num $\times$ num]( $\text{idx}(\text{zip } xs \text{ } ys) i$ ,  $\eta[i \mapsto i]$ , [.x1])
804 = CODEGENexp[ $\ell$ .num $\times$ num]( $\text{zip } xs \text{ } ys$ ,  $\eta[i \mapsto i]$ , [i, .x1])
805 = CODEGENexp[ $\ell$ .num]( $xs$ ,  $\eta[i \mapsto i]$ , [i])
806 =  $xs[i]$ 

```

Putting everything together, we get the following translation of the original **parfor** loop, which has had all the data layout combinators translated away.

```

807 #pragma omp parallel for
808 for(int i = 0; i <  $\ell$ ; i+=1) { out[i] =  $xs[i] * ys[i]$ ; }

```

A similar translation was recently presented in a more informal style by Steuwer et al. [2017]. Their experimental results show that in practice it is important to keep the indices concise and short for generating efficient imperative code and discuss how to simplify index expressions making use of range information of the indices involved.

## 5 CORRECTNESS OF THE TRANSLATION

We now justify the translation process described in Section 4.1 and Section 4.2. We have not yet formalised a correctness proof for the final translation to C (Section 4.3); this is future work.

As we stated in Section 4, the goal of the translation process was to generate a “purely imperative” command  $\mathcal{A}(\mathcal{E})_\delta(A)$  that is equivalent to the assignment command  $A :=_\delta E$ . To make this statement formal, we must define equivalence of DPIA programs. We do this in Section 5.1. We then state the functional coincidence property that establishes that our correctness property means just what we intend (Section 5.2). Finally, we prove the correctness of our translation in Section 5.3.

### 5.1 Semantics and Observational Equivalence for DPIA

We use Reddy’s relationally parametric automata-theoretic semantics of SCIR [Reddy 2013]. In this model, the interpretations of phrase types are parameterised by automata describing permitted state transitions. The model supports relationally parametric reasoning [Reynolds 1983], which enables

reasoning about locality and information hiding (following [O’Hearn and Tennent 1995]). The use of automata permits reasoning about phrases that do not affect the state (i.e. passive phrases).

Reddy’s model does not include indexed types or compound data types, as we do in DPIA, but these are straightforward to add. Indexed types are interpreted by set-theoretic functions whose codomain depends on the input (we do not interpret the data type polymorphism using parametricity because we are only interested in parametric reasoning for local state). Compound data types are interpreted as the following sets when they appear in expressions:

$$\llbracket \text{num} \rrbracket = \text{num} \quad \llbracket \delta_1 \times \delta_2 \rrbracket = \llbracket \delta_1 \rrbracket \times \llbracket \delta_2 \rrbracket \quad \llbracket \underline{\ell}.\delta \rrbracket = \{0, \dots, \ell - 1\} \rightarrow \llbracket \delta \rrbracket$$

where the set `num` is some set of number-like objects used for scalar values. This interpretation of data types allows us to straightforwardly interpret the functional primitives of Figure 3a in the model, using the standard interpretation of **map** and **reduce** explicitly given in [Steuwer et al. 2015]. This yields a coincidence property that we state in Section 5.2 below. Acceptors for compound data types are interpreted using the separating product construction, which ensures that disjoint components of compound acceptors are always non-interfering. This allows us to interpret **parfor**  $n$  as  $n$  parallel transformations on  $n$  pieces of disjoint state.

Reddy’s semantics assigns an interpretation to every DPIA phrase. For observational equivalence of DPIA programs, we are interested in closed programs. A closed program is a command phrase whose free identifiers are all have types of the form  $\text{var}[\delta]$  for possibly different  $\delta$ . Our notion of observational equivalence is standard. Two well-typed phrases  $\Delta \mid \Pi; \Gamma \vdash P, Q : \theta$  are equivalent iff for all closing contexts  $C[-]$ , the programs  $C[P]$  and  $C[Q]$ , when instantiated with the standard interpretation of variables ([Reddy 2013], Figure 2), describe the same mapping of initial to final values. We formally write  $\Delta \mid \Pi; \Gamma \vdash P \simeq Q : \theta$ , or informally  $P \simeq Q$ . Note that this relation is automatically an equivalence and is congruent with all the constructs of DPIA.

For the purposes of compilation, this notion of equivalence is justifiable: we are only interested in the relationship between the initial and final values of each variable, not the intermediate states.

## 5.2 Functional Coincidence

Our correctness criterion will have no force if we do not first establish that the assignment  $\text{out} :=_{\delta} E$  means what we think it means. We state our coincidence property formally as follows. Let  $\cdot \mid x_1 : \text{exp}[\delta_1], \dots, x_n : \text{exp}[\delta_n]; \cdot \vdash E : \text{exp}[\delta]$  be some expression phrase of DPIA built from the functional primitives in Figure 3a. Let  $\llbracket E \rrbracket : \llbracket \delta_1 \rrbracket \times \dots \times \llbracket \delta_n \rrbracket \rightarrow \llbracket \delta \rrbracket$  be the functional reference semantics of  $E$ . Use  $E$  to construct a closed phrase:

$$\cdot \mid ; v_1 : \text{var}[\delta_1], \dots, v_n : \text{var}[\delta_n], \text{out} : \text{var}[\delta] \vdash \text{out}.1 :=_{\delta} E[v_1.2/x_1, \dots, v_n.2/x_n] : \text{comm}$$

Then for all  $a_1 \in \llbracket \delta_1 \rrbracket, \dots, a_n \in \llbracket \delta_n \rrbracket, a \in \llbracket \delta \rrbracket$ , the interpretation of this command maps the store  $(v_1 \mapsto a_1, \dots, v_n \mapsto a_n, \text{out} \mapsto a)$  to the store  $(v_1 \mapsto a_1, \dots, v_n \mapsto a_n, \text{out} \mapsto \llbracket E \rrbracket(a_1, \dots, a_n))$ . In other words, this program updates the variable `out` with the result of the expression, and leaves every other variable unaffected. We use Steuwer et al. [2015]’s interpretation of the functional primitives in our DPIA semantics, so this property is immediate.

## 5.3 Correctness of the Translation from Functional to Imperative

We structure our proof by first stating a collection of equivalences that can be proved in Reddy’s model, and then use them to prove that the translation of Section 4.1 and Section 4.2 is correct (Theorem 5.1). The properties of contextual equivalences that we use in our proof, in addition to the fact that  $\simeq$  is a congruent equivalence relation, are as follows.

(1)  $\beta\eta$ -equality for non-dependent and dependent functions:

$$(\lambda x. P) Q \simeq P[Q/x] \quad P \simeq (\lambda x. P x) \quad (\Lambda x. P) \tau \simeq P[\tau/x] \quad P \simeq (\Lambda x. P x)$$

Full  $\beta\eta$ -equality for functions is one of the defining features of Idealised Algol and its descendants [Reynolds 1997]. These are all justified by Reddy's model.

(2) The **parfor** implementation of **mapParl** (Section 4.2) satisfies the following equivalence:

$$\mathbf{mapParl} N \delta_1 \delta_2 (\lambda x o. o :=_{\delta_2} F x) E A \simeq A :=_{N, \delta_2} \mathbf{mapPar} N \delta_1 \delta_2 F E \quad (5)$$

By the definition of array assignment given in Section 4.1, this property is equivalent to:

$$\mathbf{mapParl} N \delta_1 \delta_2 (\lambda x o. o :=_{\delta_2} F x) E A \simeq \mathbf{mapParl} N \delta_2 \delta_2 (\lambda x o. o :=_{\delta_2} x) (\mathbf{mapPar} \delta_1 \delta_2 F E) A$$

Expanding the definition of **mapParl**, and  $\beta$ -reducing, we must show:

$$\mathbf{parfor} N A (\lambda i o. o :=_{\delta_2} F (\mathbf{idx} N \delta_1 E i)) \simeq \mathbf{parfor} N A (\lambda i o. o :=_{\delta_2} \mathbf{idx} N \delta_2 (\mathbf{mapPar} N \delta_1 \delta_2 F E) i)$$

which is immediate from the interpretation of arrays as functions from indices to values.

(3) Reddy's model validates the following equivalence involving the use of temporary storage. For all expressions  $E$  and continuations  $C$  that are non-interfering, we have:

$$\mathbf{new} \delta (\lambda tmp. tmp.1 :=_{\delta} E; C(tmp.2)) \simeq C(E) \quad (6)$$

This equivalence relies crucially on the fact that  $C$  and  $E$  cannot interfere, so we can take a complete copy of  $E$  before invoking  $C$ . If  $C$  were able to write to storage read by  $E$ , then it would be unsafe to cache  $E$  before invoking  $C$ . In Reddy's model, we use parametricity to relate the two uses of  $C$ : one in a store that contains the state  $E$  reads, and one in a store that contains the result of evaluating  $E$ . Using parametricity and restriction to only the identity state transition on  $E$ 's portion of the store further ensures that  $C$  does not interfere with  $E$ .

(4) The **for**-loop based implementation of **reduceSeqI** should satisfy the following equivalence.

$$\mathbf{reduceSeqI} N \delta_1 \delta_2 (\lambda x y o. o :=_{\delta_2} F x y) I E C \simeq C(\mathbf{reduceSeq} N \delta_1 \delta_2 F I E) \quad (7)$$

Expanding **reduceSeqI**, and  $\beta$ -reducing, this is equivalent to showing:

$$\mathbf{new} \delta_2 (\lambda v. v.1 :=_{\delta_2} I; \mathbf{for} N (\lambda i. v.1 :=_{\delta_2} F v.2 (\mathbf{idx} E i)); C(v.2)) \simeq C(\mathbf{reduceSeq} N \delta_1 \delta_2 F I E)$$

Because the acceptor-expression pair  $v$  has been freshly allocated, it acts like a so-called "good variable" in Idealised Algol terminology. This means that the following equivalence holds, using the fact that neither  $F$  nor  $E$  interfere with  $v$ :

$$\mathbf{for} N (\lambda i. v.1 :=_{\delta_2} F v.2 (\mathbf{idx} E i)) \simeq v.1 :=_{\delta_2} \mathbf{reduceSeq} N \delta_1 \delta_2 F v.2 E$$

Now Equation 7 follows from Equation 6.

(5) Finally, we need agreement between data layout combinators and their acceptor counterparts:

$$\begin{aligned} A :=_{\delta_1 \times \delta_2} \mathbf{pair} \delta_1 \delta_2 E_1 E_2 &\simeq (\mathbf{pairAcc}_1 \delta_1 \delta_2 A :=_{\delta_1} E_1; \mathbf{pairAcc}_2 \delta_1 \delta_2 A :=_{\delta_2} E_2) \\ A :=_{N.(\delta_1 \times \delta_2)} \mathbf{zip} N \delta_1 \delta_2 E_1 E_2 &\simeq (\mathbf{zipAcc}_1 N \delta_1 \delta_1 A :=_{N, \delta_1} E_1; \mathbf{zipAcc}_2 N \delta_1 \delta_1 A :=_{N, \delta_2} E_2) \\ A :=_{N.M, \delta} \mathbf{split} N M \delta E &\simeq \mathbf{splitAcc} N M \delta A :=_{NM, \delta} E \\ A :=_{NM, \delta} \mathbf{join} N M \delta E &\simeq \mathbf{joinAcc} N M \delta A :=_{N.M, \delta} E \end{aligned}$$

The first equivalence follows directly from the definition of assignment at pair type, and  $\beta$ -reduction for pairs. The others are all straightforwardly justified in Reddy's model, given the above interpretation of acceptors for compound data types using separating products.

**THEOREM 5.1.** *The translations  $\mathcal{A}(\cdot)_\delta$  and  $C(\cdot)_\delta$  defined in Figure 4a and Figure 4b satisfy the following observational equivalences for all acceptors  $A$  and functional expressions  $E$  with disjoint sets of active identifiers:*

$$\mathcal{A}(E)_\delta(A) \simeq A :=_{\delta} E \quad C(E)_\delta(C) \simeq C(E)$$

**PROOF.** By mutual induction on the steps of the translation process. The cases for variables in both translations are immediate. The cases for the first-order combinators on numbers follow from the induction hypotheses and  $\beta$ -reduction. For example, for **negate**:

$$\mathcal{A}(\mathbf{negate} E)_{\text{num}}(A) = C(E)_{\text{num}}(\lambda x. A := \mathbf{negate} x) \simeq (\lambda x. A := \mathbf{negate} x)(E) \simeq A := \mathbf{negate} E$$

$$\frac{\ell \in \{2, 4, 8, 16\}}{\Delta \vdash \underline{\ell} : \text{width}} \qquad \frac{\Delta \vdash W : \text{width}}{\Delta \vdash \langle W \rangle \text{num} : \text{data}}$$

Fig. 6. Extension of DPIA Types with a Vector Data Type

The case for the first-order combinators in the continuation-passing translation are similar.

The acceptor-passing translation of **mapPar** uses the induction hypothesis to establish the correctness of the translations of the subterms,  $\beta$ -equality, and then the correctness property of **mapParl** (Equation 5):

$$\begin{aligned} \mathcal{A}(\text{mapPar } N \delta_1 \delta_2 F E)_{N, \delta_2}(A) &= C(\mathcal{E})_{N, \delta_1}(\lambda x. \text{mapParl } N \delta_1 \delta_2 (\lambda x o. \mathcal{A}(F x)_{\delta_2}(o)) x A) \\ &\approx \text{mapParl } N \delta_1 \delta_2 (\lambda x o. \mathcal{A}(F x)_{\delta_2}(o)) E A \\ &\approx \text{mapParl } N \delta_1 \delta_2 (\lambda x o. o :=_{\delta_2} F x) E A \\ &\approx A :=_{N, \delta_2} \text{mapPar } N \delta_1 \delta_2 F E \end{aligned}$$

The continuation-passing translation of **mapPar** relies on the acceptor-passing translation and additionally Equation 6 (temporary storage is unobservable). The acceptor-passing and continuation-passing translations for **reduceSeq** both rely on Equation 7 establishing correctness of **reduceSeqI**.

The acceptor-passing translations of the data layout combinators rely on the corresponding properties for **zip**, **split**, **join** and **pair**. The acceptor-passing cases for **fst** and **snd** follow from the induction hypothesis and  $\beta$ -equality. The correctness of the continuation-passing translations for the data layout combinators also follow by applying the induction hypothesis and using  $\beta$ -equality.  $\square$

## 6 EXTENSIONS

This section introduces two extensions to our core DPIA language: *vectorisation* and support for performing *iterative computations*. These extensions show the extensibility of DPIA for expressing more applications as well as execution strategies.

### 6.1 Vectorisation

Modern processors provide special SIMD vector instructions to perform computations on multiple scalar values in a single clock cycle. These instructions can significantly increase the performance of programs. To be able to explicitly encode strategies for vector operations in DPIA we extend the type system (Figure 6), introduce new primitives to operate on vectors (Figure 7), and extend the acceptor-passing, CPS, and code generation translations (Figure 8 and Figure 9).

We extend the definition of data types by adding a new vector type  $\langle w \rangle \text{num}$  (Figure 6) where  $w$  denotes vector width, *i.e.*, the number of elements grouped in the vector which can only be one of the sizes for which hardware instructions are available (here we assume 2, 4, 8, or 16). Widths have kind *width*, which is a subkind of *nat*. Thus a width  $w$  can be used as part of an array index.

To operate on vector values we introduce vector-specific primitives which resemble corresponding primitives defined over arrays. The **mapVec** primitive maps a scalar function over each component of a vector. Similarly to the **split** and **join** primitives, **asVector** and **asScalar** transform arrays of scalar values into arrays of vectors and vice versa. The **idxVec** primitive is used to access the individual components of a vector.

On the imperative side we add a vectorised for loop with the **forVec** primitive as well as the acceptor version of the aforementioned data layout primitives. Finally, we added the intermediate macro **mapVecI** which is expanded in our translation as follows:

$$\text{mapVecI} = \Delta n. \lambda f \text{ xs out}. \text{forVec } n \text{ out } (\lambda i o. f (\text{idxVec } n \text{ xs } i) o)$$

981 **mapVec** :  $(w : \text{width}) \rightarrow (\text{exp}[\text{num}] \rightarrow \text{exp}[\text{num}]) \rightarrow \text{exp}[\langle w \rangle \text{num}] \rightarrow \text{exp}[\langle w \rangle \text{num}]$   
 982 **asVector** :  $(w : \text{width}) \rightarrow (n : \text{nat}) \rightarrow \text{exp}[n \cdot w \cdot \text{num}] \rightarrow \text{exp}[n \cdot \langle w \rangle \text{num}]$   
 983 **asScalar** :  $(w : \text{width}) \rightarrow (n : \text{nat}) \rightarrow \text{exp}[n \cdot \langle w \rangle \text{num}] \rightarrow \text{exp}[n \cdot w \cdot \text{num}]$   
 984 **idxVec** :  $(w : \text{width}) \rightarrow \text{exp}[\langle w \rangle \text{num}] \rightarrow \text{exp}[\text{idx}[w]] \rightarrow \text{exp}[\text{num}]$   
 985 **forVec** :  $(w : \text{width}) \rightarrow \text{acc}[\langle w \rangle \text{num}] \rightarrow (\text{exp}[\text{idx}[w]] \rightarrow \text{acc}[\text{num}] \rightarrow_p \text{comm}) \rightarrow \text{comm}$   
 986 **asVectorAcc** :  $(w : \text{width}) \rightarrow (n : \text{nat}) \rightarrow \text{acc}[n \cdot \langle w \rangle \text{num}] \rightarrow \text{acc}[n \cdot w \cdot \text{num}]$   
 987 **asScalarAcc** :  $(w : \text{width}) \rightarrow (n : \text{nat}) \rightarrow \text{acc}[n \cdot w \cdot \text{num}] \rightarrow \text{acc}[n \cdot \langle w \rangle \text{num}]$   
 988 **idxVecAcc** :  $(w : \text{width}) \rightarrow \text{acc}[\langle w \rangle \text{num}] \rightarrow \text{exp}[\text{idx}[w]] \rightarrow \text{acc}[\text{num}]$   
 989 **mapVec1** :  $(w : \text{width}) \rightarrow (\text{exp}[\text{num}] \rightarrow \text{acc}[\text{num}] \rightarrow_p \text{comm}) \rightarrow \text{exp}[\langle w \rangle \text{num}] \rightarrow \text{acc}[\langle w \rangle \text{num}] \rightarrow \text{comm}$   
 990  
 991

Fig. 7. Primitives for Vectorisation

992  
 993  $\mathcal{A}(\text{mapVec } W F E)_{\langle W \rangle \text{num}}(A) = C(E)_{\langle W \rangle \text{num}}(\lambda x. \text{mapVec1 } W (\lambda x \ o. \mathcal{A}(F x)_{\text{num}}(o)) x A)$   
 994  $\mathcal{A}(\text{asVector } W N E)_{N \cdot \langle W \rangle \text{num}}(A) = \mathcal{A}(E)_{N \cdot w \cdot \text{num}}(\text{asVectorAcc } W N A)$   
 995  $\mathcal{A}(\text{asScalar } W N E)_{N \cdot w \cdot \text{num}}(A) = \mathcal{A}(E)_{N \cdot \langle W \rangle \text{num}}(\text{asScalarAcc } W N A)$   
 996  
 997  $C(\text{mapVec } W F E)_{\langle W \rangle \text{num}}(C) = \text{new } (\langle W \rangle \text{num}) (\lambda tmp. \mathcal{A}(\text{mapVec } W F E)_{\langle W \rangle \text{num}}(tmp.2); C(tmp.1))$   
 998  $C(\text{asVector } W N E)_{N \cdot \langle W \rangle \text{num}}(C) = C(E)_{N \cdot w \cdot \text{num}}(\lambda x. C(\text{asVector } W N x))$   
 999  $C(\text{asScalar } W N E)_{N \cdot w \cdot \text{num}}(C) = C(E)_{N \cdot \langle W \rangle \text{num}}(\lambda x. C(\text{asScalar } W N x))$   
 1000  
 1001

Fig. 8. Acceptor and Continuation-passing Translations for Vectorisation

1002  
 1003  
 1004  $\text{CODEGEN}_{\text{comm}}(\text{forVec } N A (\lambda i \ o. P) E, \eta) = \#pragma \text{ omp simd}$   
 1005  $\text{for}(\text{int } i = 0; i < n; i += 1) \{$   
 1006  $\text{CODEGEN}_{\text{comm}}(P[\text{idxVecAcc } N A i/o], \eta[i \mapsto i]) \}$   
 1007  
 1008  $\text{CODEGEN}_{\text{acc}[W N \cdot \text{num}]}(\text{asVectorAcc } W N A, \eta, i :: ps) = \text{CODEGEN}_{\text{acc}[N \cdot \langle W \rangle \text{num}]}(A, \eta, i/W :: ps)$   
 1009  $\text{CODEGEN}_{\text{acc}[N \cdot \langle W \rangle \text{num}]}(\text{asScalarAcc } W N A, \eta, i :: j :: ps) = \text{CODEGEN}_{\text{acc}[N \cdot w \cdot \text{num}]}(A, \eta, i * N + j :: ps)$   
 1010  $\text{CODEGEN}_{\text{acc}[\text{num}]}(\text{idxVecAcc } W A I, \eta, ps) = \text{CODEGEN}_{\text{acc}[\langle W \rangle \text{num}]}(A, \eta,$   
 1011  $\text{CODEGEN}_{\text{exp}[\text{idx}[W]]}(I, \eta, []) :: ps)$   
 1012  
 1013  $\text{CODEGEN}_{\text{exp}[M \cdot \langle N \rangle \text{num}]}(\text{asVector } N M E, \eta, i :: j :: ps) = \text{CODEGEN}_{\text{exp}[M \cdot N \cdot \text{num}]}(E, \eta, i * N + j :: ps)$   
 1014  $\text{CODEGEN}_{\text{exp}[N \cdot w \cdot \text{num}]}(\text{asScalar } W N E, \eta, i :: ps) = \text{CODEGEN}_{\text{exp}[N \cdot \langle W \rangle \text{num}]}(E, \eta, i/W :: i \% W :: ps)$   
 1015  $\text{CODEGEN}_{\text{exp}[\text{num}]}(\text{idxVec } N E I, \eta, ps) = \text{CODEGEN}_{\text{exp}[\langle W \rangle \text{num}]}(E, \eta,$   
 1016  $\text{CODEGEN}_{\text{exp}[\text{idx}[W]]}(I, \eta, []) :: ps)$   
 1017

Fig. 9. Translation of Imperative Vector Primitives DPIA to Parallel OpenMP C

1018  
 1019  
 1020 Figure 8 extends the acceptor and continuation-passing translations to account for the vector
 1021 primitives. The translations are similar to those of **mapPar**, **split**, **join**, **splitAcc**, and **joinAcc**.

1022 Code generation for vectors is defined in Figure 9. The **forVec** primitive is translated into a for
 1023 loop annotated with OpenMP's `simd` pragma which indicates to the underlying C compiler that the
 1024 instructions in the loop should be vectorised. Code generation for the other primitives is similar to
 1025 code generation for the corresponding array primitives.

1026 The assignment operation is extended to support vectors as follows:

1027  $A :=_{\langle W \rangle \text{num}} E = \text{mapVec1 } W (\lambda x \ a. a :=_{\text{num}} x) E A$   
 1028  
 1029

```

1030 iterate : (n m k : nat) → (t : data) →
1031          ((l : nat) → exp[l.n.t] → exp[l.t]) → exp[nkm.t] → exp[m.t]
1032 newDoubleBuffer : (n m k : nat) → (t : data) → exp[n.t] → acc[m.t] →
1033          (var[k.t] → comm → comm → comm) → comm
1034 truncExp : (n m : nat) → (t : data) → exp[n.t] → exp[m.t]
1035 truncAcc : (n m : nat) → (t : data) → acc[n.t] → acc[m.t]
1036 iteratelAcc : (n m k : nat) → (t : data) → acc[m.t] →
1037          ((l : nat) → acc[n.t] → exp[l.n.t] → comm) → exp[nkm.t] → comm
1038
1039
1040
1041

```

Fig. 10. Primitives for Iterative Computations

The vector extension allows us to specify strategies for which OpenMP code with vector instructions is generated. An example with the vectorised dot-product is given in Section 2 (Equation 3) and its translation into the OpenMP code is shown in Listing 2.2.

## 6.2 Iterative Computations

Our high-level functional language is deliberately restrictive. For example, it does not allow us to define recursive functions. While this severely limits the expressiveness of our language, the restrictiveness allows the generation of high performance code.

Here we show how to extend the expressiveness of the language by adding a limited form of iterative computations for which we provide an efficient implementation using a double buffering scheme. We believe that other similar extensions can be designed to extend the expressiveness without losing the ability to generate high performance code.

Figure 10 introduces the **iterate** primitive which applies a given function a fixed number of times to an array, by supplying the output of one function application as input to the next. We take our formulation of **iterate** from Steuwer et al. [2015] with a type which allows the size of the array to shrink by a fixed factor  $n$  in each iteration, allowing **iterate** to be used for implementing a parallel tree-based reduction by iteratively performing multiple sequential reductions in parallel – an implementation strategy suggested by Nvidia to be used on their GPUs [Harris 2008].

We implement **iterate** by double buffering, that is, we allocate two intermediate buffers and switch between which is read from and written to after each iteration. For this we introduce the **newDoubleBuffer** primitive:

$$\mathbf{newDoubleBuffer} \ N \ M \ K \ \delta \ E \ A \ (\lambda v \ \mathit{swap} \ \mathit{done} . P)$$

Similar to the **new** primitive it binds a variable  $v$  which can be read from and written to, but in contrast to **new** this variable provides an interface to two different buffers in memory: reading from one and writing to the other. Besides the variable  $v$  the primitive also takes two commands: *swap* and *done*. The intention is that when *swap* is invoked, it swaps which buffer is read from and which buffer is written to and when *done* is invoked, the reading buffer is swapped as usual but the writing buffer is set directly to the output.

Figure 12 shows the code generation for **newDoubleBuffer** showing its efficient implementation. After allocation of two temporary buffers, two pointers are initialised: one pointing to the input  $E$  and the other to the first temporary buffer `tmp1`. The implementation of *swap* and *done* are added to an extended environment  $\eta$  which as well as mapping from DPIA identifiers to C variable names, now also maps identifiers of type `comm` to C statements. Every time *swap* or *done* are used inside of  $P$ , the statements from the environment will be emitted. We implement *swap* efficiently by swapping the two pointers between the two temporary buffers; a flag is toggled every time to

1079  $\mathcal{A}(\mathbf{iterate} \ N \ M \ K \ \delta \ F \ E)_{N.\delta}(A) = C(\llbracket E \rrbracket_{N^k N.\delta}(\lambda x. \mathbf{iterateIAcc} \ N \ M \ K \ \delta \ A \ (\lambda l. \lambda o \ x. \mathcal{A}(\llbracket F \ I \ x \rrbracket_{l.\delta}(o)) \ x))$   
 1080  $C(\mathbf{iterate} \ N \ M \ K \ \delta \ F \ E)_{N.\delta}(C) = \mathbf{new} \ (N.\delta) \ (\lambda tmp. \mathcal{A}(\mathbf{iterate} \ N \ M \ K \ \delta \ F \ E)_{N.\delta}(tmp.2); C(tmp.1))$   
 1081

1082 Fig. 11. Acceptor and Continuation-passing Translations for the **iterate** primitive  
 1083

1084  $\text{CODEGEN}_{\text{comm}}(\mathbf{newDoubleBuffer} \ N \ M \ K \ \delta \ E \ A \ (\lambda p. P), \eta) =$   
 1085  $\{$   
 1086  $\text{CODEGEN}_{\text{data}}(K.\delta) \ tmp1; \ \text{CODEGEN}_{\text{data}}(K.\delta) \ tmp2;$   
 1087  $\text{CODEGEN}_{\text{data}}(\delta) * \text{in\_ptr} = \&(\text{CODEGEN}_{\text{exp}}[N.\delta](E, \eta, [0]));$   
 1088  $\text{CODEGEN}_{\text{data}}(\delta) * \text{out\_ptr} = \text{tmp1};$   
 1089  $\text{char} \ \text{flag} = 1;$   
 1090  $\text{CODEGEN}_{\text{comm}}(P[\langle\langle\langle\langle (v_a, v_e), \text{swap} \rangle, \text{done} \rangle\rangle\rangle/p],$   
 1091  $\eta[v_a \mapsto \text{out\_ptr}, v_e \mapsto \text{in\_ptr},$   
 1092  $\text{swap} \mapsto \{\text{in\_ptr} = \text{flag} ? \text{tmp1} : \text{tmp2};$   
 1093  $\text{out\_ptr} = \text{flag} ? \text{tmp2} : \text{tmp1};$   
 1094  $\text{flag} = \text{flag} \wedge 1;\},$   
 1095  $\text{done} \mapsto \{\text{in\_ptr} = \text{flag} ? \text{tmp1} : \text{tmp2};$   
 1096  $\text{out\_ptr} = \&(\text{CODEGEN}_{\text{acc}}[M.\delta](A, \eta, [0])); \})$   
 1097  $\}$

1098 Fig. 12. Translation of the Imperative **newDoubleBuffer** primitive to Parallel OpenMP C  
 1099

1100 indicate how to configure the buffers. We implement *done* by setting the pointer used for writing  
 1101 to the output acceptor *A*.

1102 The **iterate** primitive is implemented in terms of **newDoubleBuffer** via a translation to the  
 1103 **iterateIAcc** macro given in Figure 11. **iterateIAcc** is then defined as follows:  
 1104

1105  $\mathbf{iterateIAcc} = \lambda n \ m \ k \ \delta \ (\lambda a \ f \ e. \mathbf{newDoubleBuffer} \ (n^k m) \ m \ (n^k m) \ \delta \ e \ a \ (\lambda v \ \text{swap} \ \text{done}.$   
 1106  $\mathbf{for} \ k \ (\lambda i. f \ (n^{k-1} m) \ (\mathbf{truncAcc} \ (n^k m.\delta) \ (n^{k-i-1} m.\delta) \ \delta \ v.1)$   
 1107  $\ (\mathbf{truncExp} \ (n^k m.\delta) \ (n^{k-i} m.\delta) \ \delta \ v.2);$   
 1108  $\mathbf{if} \ (i < k - 1) \ \mathbf{then} \ \text{swap} \ \mathbf{else} \ \text{done}))$

1109 We assume a standard extension to support booleans. The **for** loop performs the iterative application  
 1110 of the function *F*. After each iteration *swap* is invoked until the second to last iteration when *done*  
 1111 is invoked. This implementation ensures that no unnecessary copies of input or output arrays are  
 1112 made: the first iteration reads straight from the input and the last writes directly to the output.

1113 The size of the arrays change in each iteration but only two buffers of a fixed size are allocated. To  
 1114 ensure that the program typechecks we introduce **truncExp** and **truncAcc** primitives, which are  
 1115 no-ops that restrict the size of a larger array. For the implementation of **iterateIAcc** we also extend  
 1116 the arithmetic expressions with a subtraction operation. All nat expressions in the implementation  
 1117 are guaranteed to be non-negative as *i* is the iteration variable and thus always smaller than *k*.  
 1118

## 1119 7 EXPERIMENTAL RESULTS

1120 We implemented DPIA and the compilation to OpenMP code via the translation described in  
 1121 Section 4. This section evaluates the quality of the OpenMP code generated with our implementation.  
 1122 We are interested to see if this formal translation introduces overheads compared to hand-written  
 1123 OpenMP code. We start by describing our experimental setup and the benchmarks used.

1124 *Experimental Setup.* We used an Intel Core i7-7700 CPU with 4 physical cores and hyper-threading  
 1125 enabled. We measured the runtime of 10 runs and report the median runtime. We compare against  
 1126

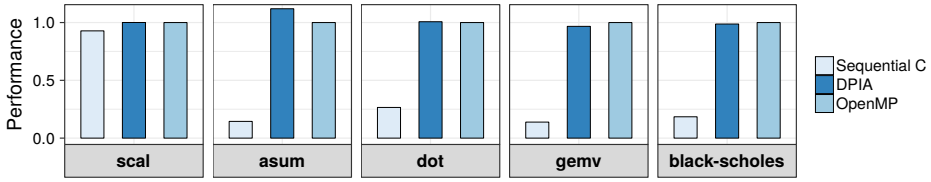


Fig. 13. Performance of DPIA generated code for five benchmarks compared to sequential C and OpenMP.

hand-written OpenMP code on four simple linear algebra benchmarks as well as the Black-Scholes application from financial stochastics.

*Performance Results.* Figure 13 shows the relative performance of sequential C code on the left, the code generated from out DPIA language in the middle and hand-written OpenMP code on the right. As can be seen, parallelism is beneficial for most of our benchmarks and our implementation does not introduce any significant overhead compared to the hand-written OpenMP code. For the absolute sum benchmark (asum) our DPIA implementation even slightly outperforms the hand-written OpenMP implementation as we use a more complex strategy using vectorisation along the lines of the example shown in Section 2.

## 8 RELATED WORK

*Idealised Algol and Syntactic Control of Interference.* We have heavily relied on Reynolds’ insightful design of Idealised Algol (IA) in this work, originally spelt out in [Reynolds 1997]. IA’s orthogonal combination of typed  $\lambda$ -calculus and imperative programming has given us the ideal language in which to formalise compilation from functional to imperative code. Moreover, Reynolds’ Syntactic Control of Interference (SCI) [Reynolds 1978] enabled us to ensure that we always produce deterministic data race free programs. Brookes describes a concurrent version of Idealised Algol [Brookes 2002] that he calls “parallel”, but is intentionally a non-deterministic concurrent language that allows threads to communicate through shared memory.

Reynolds presents SCI as a series of principles for a language design, which are formulated as a substructural type system by O’Hearn et al. [1999]. We have used O’Hearn et al.’s formulation in our design of DPIA. We do not know of any other work using IA or SCI as an intermediate language for compilation, although Ghica has used interference controlled Idealised Algol as a high-level language for hardware synthesis [Ghica 2007].

Reynolds’ original presentation of IA describes its semantics in terms of a functor categories [Oles 1982]. This semantics models the stack discipline of IA (which Reynolds uses to systematically derive a compilation strategy for IA [Reynolds 1995]), but does not model non-interference or the locality of fresh state, both of which we rely upon in Section 5. O’Hearn and Tennent were the first to use relational parametricity to model locality [O’Hearn and Tennent 1995]. Models of non-interference are given by O’Hearn [1993], Tennent [1990], and, O’Hearn and Tennent [1993], refining the original Reynolds/Oles functor category semantics. The semantic insights in this work later fed into Separation Logic [Ishtiaq and O’Hearn 2001]. Reddy’s object-space semantics [Reddy 1994, 1996] also modelled non-interference, but took an intensional viewpoint based on modelling interactions with the store, rather than transformations of the store. Reddy later synthesised the relationally parametric and intensional approaches in his automata-theoretic model [Reddy 2013], which we used in Section 5.



1177 *Functional Compilation Approaches Targeting Parallel Architectures*. There exist multiple func-  
1178 tional approaches for generating code for parallel hardware. Steuwer et al. [2015] use a data parallel  
1179 language similar to the functional subset of DPIA. Semantics preserving rewrite rules are used to  
1180 explore the space of possible implementations showing that achieving high performance across  
1181 multiple architectures from functional code is possible. Obsidian [Svensson et al. 2016] is a func-  
1182 tional low-level GPU programming language which gives programmers flexibility over how to  
1183 write efficient code while still providing functional abstractions. Accelerate is a Haskell embedded  
1184 domain specific language providing higher level abstractions with the aim of generating efficient  
1185 parallel code [Chakravarty et al. 2011; McDonell et al. 2013]. LiquidMetal [Dubach et al. 2012]  
1186 targets Field-Programmable Gate Array (FPGAs) and GPUs by extending Java with a data-flow  
1187 programming model with built-in functional map and reduce operators. Bergstrom and Reppy  
1188 [2012] compile NESL, which is a first-order dialect of ML supporting nested data-parallelism and  
1189 introduced by Blelloch [1993], to GPU code. Anderson et al. [2017] present an approach to compile  
1190 array-code written in Julia to parallel code. Delite [Sujeeth et al. 2014] enables the creation of  
1191 domain-specific languages using functional parallel patterns targeting parallel hardware.

1192

## 1193 9 CONCLUSIONS AND FUTURE WORK

1194 We introduced a novel data parallel programming language DPIA, an interference controlled dialect  
1195 of Idealised Algol. We showed how DPIA is used as a foundation for a formal translation of data  
1196 parallel functional programs to imperative code for parallel machines. The approach offers strong  
1197 guarantees about the absence of data-races in the generated programs and a straightforward  
1198 translation strategy for generating efficient OpenMP code. Our experimental results show that our  
1199 approach yields efficient code on a par with hand-coded low-level parallel code.

1200 While this paper has presented an implementation targeting OpenMP we are eager to extend this  
1201 work to target a diverse set of heterogeneous devices such as GPUs and FPGAs where the choice  
1202 of strategy is more challenging than on traditional CPUs. We envision such an extended version  
1203 of DPIA as a unifying language for generating efficient code for diverse parallel hardware from a  
1204 single high-level portable program, thus obtaining performance portability. To pursue our vision  
1205 we intend to model parallelisation and memory hierarchies present in modern parallel devices such  
1206 as GPUs in the type system which is for example required to guarantee deadlock freedom. While  
1207 OpenMP follows the sequential-by-default fork-join parallelism model, the dominant low-level GPU  
1208 programming models OpenCL and CUDA follow a parallelism-by-default model which requires  
1209 changes to the presented translation to continue to guarantee race-freedom. Finally, we plan to  
1210 push our formalisation further to prove the correctness of the final translation from DPIA to C,  
1211 OpenMP, and in the future OpenCL.

1212

## 1213 REFERENCES

- 1214 Todd A. Anderson, Hai Liu, Lindsey Kuper, Ehsan Toton, Jan Vitek, and Tatiana Shpeisman. 2017. Parallelizing Julia  
1215 with a Non-Invasive DSL. In *31st European Conference on Object-Oriented Programming, ECOOP 2017, June 19-23, 2017,*  
1216 *Barcelona, Spain (LIPICs)*, Peter Müller (Ed.), Vol. 74. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 4:1-4:29.  
1217 <https://doi.org/10.4230/LIPICs.ECOOP.2017.4>
- 1218 Lars Bergstrom and John H. Reppy. 2012. Nested data-parallelism on the gpu. In *ICFP*. ACM, 247-258.
- 1219 Guy E. Blelloch. 1993. *NESL: A nested data-parallel language (version 2.6)*. Technical Report CMU-CS-93-129. School of  
1220 Computer Science, Carnegie Mellon University.
- 1221 Stephen D. Brookes. 2002. The Essence of Parallel Algol. *Inf. Comput.* 179, 1 (2002), 118-149.
- 1222 Bryan Catanzaro, Michael Garland, and Kurt Keutzer. 2011. Copperhead: Compiling an Embedded Data Parallel Language  
(PPoPP). ACM, 10.
- 1223 Manuel M. T. Chakravarty, Gabriele Keller, Sean Lee, Trevor L. McDonell, and Vinod Grover. 2011. Accelerating Haskell  
1224 array codes with multicore GPUs. In *DAMP*. ACM, 3-14.

1225

- 1226 Olivier Danvy, Kevin Millikin, and Lasse R. Nielsen. 2007. On one-pass CPS transformations. *J. Funct. Program.* 17, 6 (2007),  
1227 793–812.
- 1228 Christophe Dubach, Perry Cheng, Rodric Rabbah, David F. Bacon, and Stephen J. Fink. 2012. Compiling a High-level  
1229 Language for GPUs: (via Language Support for Architectures and Compilers) (*PLDI*). ACM, 12.
- 1230 G. Gentzen. 1935. Untersuchungen über das logische Schließen. I. *Mathematische zeitschrift* 39, 1 (1935).
- 1231 Dan R. Ghica. 2007. Geometry of synthesis: a structured approach to VLSI design. In *POPL*. ACM, 363–375.
- 1232 Mark Harris. 2008. *Optimizing Parallel Reduction in CUDA*. Technical Report. nVidia. <http://developer.download.nvidia.com/assets/cuda/files/reduction.pdf>
- 1233 Samin S. Ishtiaq and Peter W. O’Hearn. 2001. BI as an Assertion Language for Mutable Data Structures. In *POPL*. ACM,  
1234 14–26.
- 1235 Steve Klabnik and Carol Nichols. 2018. *The Rust Programming Language* (2nd ed.). No Starch Press.
- 1236 Trevor L. McDonnell, Manuel M. T. Chakravarty, Gabriele Keller, and Ben Lippmeier. 2013. Optimising purely functional  
1237 GPU programs. In *ICFP*. ACM, 49–60.
- 1238 Shayan Najd, Sam Lindley, Josef Svenningsson, and Philip Wadler. 2016. Everything old is new again: quoted domain-specific  
1239 languages. In *PEPM*. ACM, 25–36.
- 1240 Peter W. O’Hearn. 1993. A Model for Syntactic Control of Interference. *Mathematical Structures in Computer Science* 3, 4  
1241 (1993), 435–465.
- 1242 Peter W. O’Hearn, John Power, Makoto Takeyama, and Robert D. Tennent. 1999. Syntactic Control of Interference Revisited.  
1243 *Theor. Comput. Sci.* 228, 1-2 (1999), 211–252.
- 1244 Peter W. O’Hearn and Robert D. Tennent. 1993. Semantical Analysis of Specification Logic, 2. *Inf. Comput.* 107, 1 (1993),  
1245 25–57.
- 1246 Peter W. O’Hearn and Robert D. Tennent. 1995. Parametricity and Local Variables. *J. ACM* 42, 3 (1995), 658–709.
- 1247 Frank J. Oles. 1982. *A Category Theoretic Approach To Semantics of Programming Languages*. Ph.D. Dissertation. Syracuse  
1248 University, Syracuse, New York.
- 1249 D. Prawitz. 1965. *Natural Deduction: A Proof-Theoretical Study*. Almqvist and Wiksell.
- 1250 Jonathan Ragan-Kelley, Andrew Adams, Dillon Sharlet, Connelly Barnes, Sylvain Paris, Marc Levoy, Saman P. Amarasinghe,  
1251 and Frédo Durand. 2018. Halide: decoupling algorithms from schedules for high-performance image processing. *Commun.*  
1252 *ACM* 61, 1 (2018), 106–115. <https://doi.org/10.1145/3150211>
- 1253 Uday Reddy. 1994. Passivity and independence. In *Proceedings of the Ninth Annual IEEE Symp. on Logic in Computer Science,*  
1254 *LICS 1994*, Samson Abramsky (Ed.). IEEE Computer Society Press, 342–352.
- 1255 Uday S. Reddy. 1996. Global State Considered Unnecessary: An Introduction to Object-Based Semantics. *Lisp and Symbolic*  
1256 *Computation* 9, 1 (1996), 7–76.
- 1257 Uday S. Reddy. 2013. Automata-Theoretic Semantics of Idealized Algol with Passive Expressions. *Electr. Notes Theor. Comput.*  
1258 *Sci.* 298 (2013), 325–348.
- 1259 John C. Reynolds. 1978. Syntactic Control of Interference. In *POPL*, Alfred V. Aho, Stephen N. Zilles, and Thomas G.  
1260 Szymanski (Eds.). ACM Press, 39–46.
- 1261 J. C. Reynolds. 1983. Types, Abstraction and Parametric Polymorphism. In *Proc. IFIP 9th World Computer Congress (Information*  
1262 *Processing)*, Vol. 83. 513–523.
- 1263 John C. Reynolds. 1995. Using Functor Categories to Generate Intermediate Code. In *POPL*. ACM Press, 25–36.
- 1264 John C. Reynolds. 1997. The Essence of Algol. In *Algol-like Languages*, Peter W. O’Hearn and Robert D. Tennent (Eds.).  
1265 Birkhäuser Boston, 67–88.
- 1266 John C. Reynolds. 1998. Definitional Interpreters for Higher-Order Programming Languages. *Higher-Order and Symbolic*  
1267 *Computation* 11, 4 (1998), 363–397.
- 1268 G. L. Steele. 1978. *Rabbit: A Compiler for Scheme*. Master’s thesis. MIT.
- 1269 Michel Steuwer, Christian Fensch, Sam Lindley, and Christophe Dubach. 2015. Generating performance portable code using  
1270 rewrite rules: from high-level functional expressions to high-performance OpenCL code (*ICFP 2015*). ACM, 205–217.
- 1271 Michel Steuwer, Toomas Rimmelg, and Christophe Dubach. 2016. Matrix multiplication beyond auto-tuning: rewrite-based  
1272 GPU code generation. In *2016 International Conference on Compilers, Architectures and Synthesis for Embedded Systems,*  
1273 *CASES 2016, Pittsburgh, Pennsylvania, USA, October 1-7, 2016*. ACM, 15:1–15:10. <https://doi.org/10.1145/2968455.2968521>
- 1274 Michel Steuwer, Toomas Rimmelg, and Christophe Dubach. 2017. Lift: A Functional Data-Parallel IR for High-Performance  
GPU Code Generation (*CGO 2017*). IEEE, 74–85.
- Arvind K. Sujeeth, Kevin J. Brown, Hyoukjoong Lee, Tiark Rompf, Hassan Chafi, Martin Odersky, and Kunle Olukotun.  
2014. Delite: A Compiler Architecture for Performance-Oriented Embedded Domain-Specific Languages. *ACM TECS* 13,  
4s, Article 134 (2014), 25 pages.
- Bo Joel Svensson, Ryan R. Newton, and Mary Sheeran. 2016. A language for hierarchical data parallel design-space  
exploration on GPUs. *J. Funct. Program.* 26 (2016), e6.

- 1275 Robert D. Tennent. 1990. Semantical Analysis of Specification Logic. *Inf. Comput.* 85, 2 (1990), 135–162.
- 1276 Philip W. Trinder, Kevin Hammond, Hans-Wolfgang Loidl, and Simon L. Peyton Jones. 1998. Algorithms + Strategy =  
1277 Parallelism. *J. Funct. Program.* 8, 1 (1998), 23–60. <http://journals.cambridge.org/action/displayAbstract?aid=44147>
- 1278
- 1279
- 1280
- 1281
- 1282
- 1283
- 1284
- 1285
- 1286
- 1287
- 1288
- 1289
- 1290
- 1291
- 1292
- 1293
- 1294
- 1295
- 1296
- 1297
- 1298
- 1299
- 1300
- 1301
- 1302
- 1303
- 1304
- 1305
- 1306
- 1307
- 1308
- 1309
- 1310
- 1311
- 1312
- 1313
- 1314
- 1315
- 1316
- 1317
- 1318
- 1319
- 1320
- 1321
- 1322
- 1323