# Designing Data Layout Specification Languages

Jakob Schuster

October 2023

## 1   Introduction

Porgrammers writing software often enjoy the conceptual abstraction of types and values, without considering the actual layout of their values in memory, on a bit- and byte-level. The layout of types is typically reserved for the compiler, handled under the programmers' feet. However, some low-level programmers require control over how their objects are layed out. When writing device drivers, file systems, and network software, systems programmers need to closely design their data layout for compatibility or performance reasons. In these situations, they must typically throw away the abstraction of the data type, writing code which performs complex bit-twiddling operations with tedious and error-prone implementations. We desire a best-of-both-worlds solution; where the programmer can design the layout of their values, without giving up safety to manage the getting and setting of their structures by hand.

One solution to this is the idea of data layout specification, as a language feature or extension. In this approach, the programmer dictates to the compiler how it should lay out values of a particular type, and the compiler generates appropriate getter and setter functions reliably following their specification. DARGENT is one leading example, applying this idea to the functional language COGENT and verifying its marshalling operations ([1]). Another similar approach by Baudon et al. introduces RIBBIT, generating code targeting LLVM IR ([2]).

In this report, we contribute to the development of data layout languages and discuss their application and formalisation. First, we extend DARGENT to allow for the laying out of recursive types. Then, we propose a data layout language for C, giving a static semantics of its core syntax to define layout well-formedness and a matching relation between layouts and types. In doing this, we give an account of how C automatically generates layouts for its types. In the design and formalisation of this language, we follow the example of DARGENT, essentially reproducing the steps of Chen's formalisation in a new language setting ([3]). Finally, we discuss more broadly the relationship of layouts and types, speculating on how programming languages could treat data layouts more thoroughly to allow for greater portability and interoperability.

# 2 Preliminaries

## Notation

We assume the knowledge of basic set theory concepts, and reserve this section for specific notation that appears throughout this report.

**Ordered Lists**   We use overlines $\overline{x_i}$ to represent an ordered list of elements. Often, we treat lists like $Q = \overline{k : v}$ as *environments*, where looking up a key $k$ returns its value $v$. This lookup operation is notated as $Q.k = v$.

**Inference Rules**   Most static semantics are given through inference rules, where premises are written above the line and the conclusion is written below the line. The rules below provide an example.

$$\frac{}{0 \; \mathsf{nat}}\textsc{ZeroNat} \qquad \frac{n \; \mathsf{nat}}{\mathsf{s}\,(n) \; \mathsf{nat}}\textsc{SuccNat}$$

**Denotation Brackets**   Where we give denotational semantics, we use denotation brackets $[\![ \cdot ]\!]$ around the syntactic construct.

## Cogent and Dargent

COGENT is a functional programming language which is transpiled to C. It is designed with systems programming in mind, and features a uniqueness type system which allows for two different semantics: one imperative semantics useful in C code generation, and one purely functional semantics which simplifies the process of verifying correctness properties. These aspects of COGENT are central to its purpose, but they will not be the subject of our discussion in this report. We will consider COGENT in the limited capacity of being the host language to DARGENT, a data layout language which allows COGENT programmers to dictate to the compiler how their types should be laid out. DARGENT's design is declarative; memory specifications are written once by the programmer, with no need to write getters and setters or worry about low-level memory details anywhere else in the program. The actual implementation and management of data with custom layouts is handled by the compiler. DARGENT also fits into the COGENT verification framework, generating correctness proofs.

$$
\begin{array}{llll}
\text{Types} & \tau & ::= & T & \text{(primitive types)} \\
& & | & () & \text{(unit)} \\
& & | & \alpha & \text{(type variables)} \\
& & | & t_r & (\ \text{recursive type variables}\ ) \\
& & | & \mu\ \{\overline{f_i : \tau_i}\}\ s & (\ \text{recursive}\ \text{records)} \\
& & | & \langle \overline{A_i\ \tau_i}\rangle s & \text{(variants)} \\
& & | & \mathsf{K}\ \overline{\tau_i}\ s & \text{(abstract types)} \\
& & | & \dots & \\
\text{Primitive types} & T & ::= & \mathsf{U}n \mid \mathsf{Bool} & (\mathrm{n} = 1\dots64) \\
\text{Field names} & & \ni & \mathsf{f} & \\
\text{Constructors} & & \ni & \mathsf{A} & \\
\text{Sigils} & s & ::= & ⓑ \mid ⓤ & \\
& ⓑ & ::= & ⓦ \mid ⓡ & \\
\text{Recursive parameters} & \mu & ::= & \mathbf{mu}\ t. \mid \epsilon &
\end{array}
$$

(lists are represented by $\overline{\text{overlines}}$)

Figure 1: The syntax of COGENT types, with recursive parameters.

# 3 Extending Dargent to Recursive Types

With the addition of recursive types ($\mu$-types), records can be optionally prefixed with a recursive parameter (Figure 1). Recursive parameter variables $t_r$ can occur within $\mu$-types. The approach introduced by Murray lays out $\mu$-types using pointers, as opposed to a 'flattened' memory layout ([4]). As a result, an intuitive DARGENT layout can be written using pointers; explicit recursion need not be present in the memory layout of a $\mu$-type.

```
layout Number = record {
  l : record {
    variant (1b) { Zero(0) : 0B, Succ(1) : pointer }
  }
}
```

To allow for this new use of syntax, we modify DARGENT's main typing rule ([1]), with our modifications highlighted :

$$
\frac{\begin{array}{c} \overline{\mathsf{f}_i}\ \text{distinct} \qquad \text{for each } i\colon \Delta; C \vdash \tau_i\ \mathbf{wf} \\ \ell\ \mathbf{wf} \qquad \ell \sim \mu\ \{\overline{\mathsf{f}_i : \hat{\tau}_i}\} \qquad (\ell, \mu\ \{\overline{\mathsf{f}_i : \hat{\tau}_i}\}) \in C \end{array}}{\Delta; C \vdash \mu\ \{\overline{\mathsf{f}_i :: \tau_i}\}b_\ell\ \mathbf{wf}}\text{RecWf}
$$

Clearly, the only change is the insertion of $\mu$-parameters. We also tweak the matching relation between layouts and types. We account for the new record syntax in URECORDMATCH, and we capture the relation of $\mu$-type parameter variables $t_r$ as pointers in the MUPARAMMATCH rule. (Figure 2). Since the surface syntax of DARGENT needs no alteration, we can leave untouched the layout-wellformedness and code generation aspects of DARGENT.

$$\boxed{\ell \sim \hat{\tau}}$$

$$\frac{\text{bits}\ (T) = s}{\text{BitRange}\ (o, s) \sim T}\text{P\scriptsize RIM}\text{T\scriptsize Y}\text{M\scriptsize ATCH} \qquad \begin{aligned}\text{bits}\ (\mathsf{U}n) &= n \\ \text{bits}\ (\mathsf{Bool}) &= 1\end{aligned}$$

$$\frac{M \text{ is the pointer size}}{\text{BitRange}\ (o, M) \sim \mathsf{pointer}}\text{B\scriptsize OXED}\text{T\scriptsize Y}\text{M\scriptsize ATCH} \qquad \frac{M \text{ is the pointer size}}{\text{BitRange}\ (o, M) \sim t_r}\text{M\scriptsize U}\text{P\scriptsize ARAM}\text{M\scriptsize ATCH}$$

$$\frac{}{x\{o\} \sim \hat{\tau}}\text{V\scriptsize AR}\text{M\scriptsize ATCH}$$

$$\frac{}{()\sim ()}\text{U\scriptsize NIT}\text{M\scriptsize ATCH} \qquad \frac{\text{for each } i:\ \ell_i \sim \hat{\tau}_i}{\mathsf{record}\ \{\overline{\mathsf{f}_i : \ell_i}\} \sim \mu\ \{\overline{\mathsf{f}_i : \hat{\tau}_i}\}}\text{U\scriptsize R}\text{ECORD}\text{M\scriptsize ATCH}$$

$$\frac{\text{for each } i:\ \ell_i \sim \hat{\tau}_i}{\mathsf{variant}\ (s)\ \{\overline{\mathsf{A}_i\ (v_i) : \ell_i}\} \sim \langle \overline{\mathsf{A}_i\ \hat{\tau}_i} \rangle}\text{V\scriptsize ARIANT}\text{M\scriptsize ATCH}$$

Figure 2: Matching relation between layouts and types, extended to $\mu$-types.

# 4 Dargent for C

The idea of surface syntax to control the layout of types can be applied to languages beyond COGENT. In this section, we apply it to the widely-used systems programming language C, as it is a language in which programmers often apply non-intuitive bit-twiddling approaches to manually manage the memory of their data. We formalise a syntax and its associated static semantics to afford C programmers control over the memory layout of their types.

In our syntax, each type can be associated with a desired memory representation using the as keyword (Figure 3). The position of each field of a struct can be moved to a position in memory, either absolutely using at, or relative to another field using after/before. As shown in the example above, this positioning system allows programmers to break C's traditional padding specifications and pack structs to meet their needs. In the following sections, we formalise this language extension.

### Bits and Bytes

First, we define a generic language of bits and bytes, and operations thereupon, which we will use across our later language definitions. Memory quantities (typically representing either the size or offset of some object) are expressed as quantities of *bits* $m\mathsf{b}$ and *bytes* $n\mathsf{B}$.

$$\begin{aligned}\text{Memory } o, s &::= n\mathsf{B} \mid m\mathsf{b} \mid n\mathsf{B}\ m\mathsf{b} \\ \text{Naturals} &\ni n, m\end{aligned}$$

Conversion functions are defined between sizes and the number of bits they denote in $\mathbb{N}$, allowing us to define addition, subtraction, scalar multiplication and the modulo operation on sizes.

```
struct {
    struct { int a ; _Bool b ; } fst;
    struct { char c; }* ptr ;
    union { char d ; char e [3]; } sum ;
} as {
    fst: { a : 4B after b; b : 1B at 0B; };
    ptr : 8B at 8B;
    sum : 3B at 5B;
} example;
```
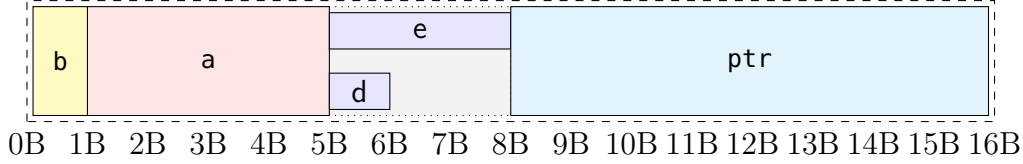
Figure 3: A demo of the C layout language.

$$
\begin{aligned}
\text{bits} \left( n\mathsf{B} \ m\mathsf{b} \right) &= \text{bits} \left( n\mathsf{B} \right) + \text{bits} \left( m\mathsf{b} \right) \\
\text{bits} \left( n\mathsf{B} \right) &= 8 \cdot n \\
\text{bits} \left( m\mathsf{b} \right) &= m
\end{aligned}
$$

$$
\text{size} \left( n \right) = \left\lfloor \tfrac{n}{8} \right\rfloor \mathsf{B} \ (n \bmod 8)\mathsf{b}
$$

$$
\begin{aligned}
s_1 + s_2 &= \text{size} \left( \text{bits} \left( s_1 \right) + \text{bits} \left( s_2 \right) \right) \\
s_1 - s_2 &= \text{size} \left( \text{bits} \left( s_1 \right) - \text{bits} \left( s_2 \right) \right) \quad [s_1 \geq s_2] \\
s \cdot n &= \text{size} \left( \text{bits} \left( s \right) \cdot n \right) \\
\tfrac{s_1}{s_2} &= \tfrac{\text{bits}(s_1)}{\text{bits}(s_2)} \\
s \bmod n &= \text{size} \left( \text{bits} \left( s \right) \bmod n \right)
\end{aligned}
$$

$$
\text{aligned} \left( o, s \right) = s \cdot \left\lceil \tfrac{o}{s} \right\rceil
$$

We define a function aligned $(o, s)$ to handle C's notion of alignment, finding the next alignment after $o$ which is a multiple of an alignment size $s$. For example, aligned $(3\mathsf{B} \ 7\mathsf{b}, 2\mathsf{B}) = 4\mathsf{B}$; three *padding bits* are required to align $3\mathsf{B} \ 7\mathsf{b}$ to the nearest $2\mathsf{B}$. We also define a *range* $(\!| \ o, s \ |\!)$ to represent a region of bits starting at offset $o$ with size $s$.

$$
\text{Ranges } r ::= (\!| \ o, s \ |\!)
$$

We also have a domain of *fields*, shared across all of the language fragments we discuss.

$$
\text{Fields } \ni f
$$

# Types in C

Next, we give the type system of C. C's surface syntax as outlined in the specification makes it difficult to separate type information from the other language features of C; in our formalisation, taking strong inspiration from CompCert ([5]), the types have been isolated for simplicity (Figure 4). One C feature worth discussing is the assignment of attributes $\alpha$, allowing the programmer to add special behaviours to variables, such as aligning one type as though it were another in memory.

| Types | $\tau$ | ::= | void | (the void type) |
|---|---|---|---|---|
| | | \| | $\mathsf{char}(\sigma)$ | (character types) |
| | | \| | $\mathsf{int}(\sigma, s_i)$ | (integer types) |
| | | \| | float | (floating-point) |
| | | \| | $\mathsf{double}(s_d)$ | (double-precision int) |
| | | \| | $\mathsf{array}\,(\tau, n)$ | (array of $\tau$) |
| | | \| | $\mathsf{struct}\left(\overline{f_i : \tau_i}\right)$ | (structure) |
| | | \| | $\mathsf{union}\left(\overline{f_i : \tau_i}\right)$ | (union) |
| | | \| | $\mathsf{pointer}\,(\tau)$ | (pointer to $\tau$) |
| | | \| | $\mathsf{enum}\left(\overline{f_i}\right)$ | (enumeration type) |
| | | \| | $\beta$ | (type variable) |
| | | \| | $\mathsf{fun}\,(\overline{\tau_i}, \overline{\tau_j}, \tau_r)$ | (function with $\overline{\tau_i}$ man, $\overline{\tau_j}$ opt args) |
| | | \| | $\mathsf{bitfield}\left(\tau, \overline{f_i : r}\right)$ | (bitfield with basic type $\tau$, fields $\overline{f_i}$, ranges $\overline{r_i}$) |
| | | \| | $\tau\,(\alpha)$ | ($\tau$ with attribute $\alpha$) |
| | | \| | $\tau$ as $l$ | ($\tau$ layed out as $l$) |
| Signedness | $\sigma$ | ::= | signed \| unsigned | |
| Int sizes | $s_i$ | ::= | bool \| short \| default | (valid integer sizes) |
| | | \| | long \| long long | |
| Double sizes | $s_d$ | ::= | default \| long | |
| Attributes | $\alpha$ | ::= | align as $\tau$ \| volatile \| $\cdots$ | (optionally, align the type as though it were $\tau$) |
| Type variables | $\ni$ | | $\beta$ | |

Figure 4: Type system of C, adapted from the CompCert formalisation, with some omissions to focus on memory-relevant properties.

## Bit-fields

Another C type worth discussing is the bit-field, a language feature with highly particular memory behaviour. A bit-field can be thought of as a region whose size is a multiple of some basic type $\tau$ (usually an integer), with specific subsections of this region given names, and others marked as empty padding (Figure 5). Special syntax :0 pads the bit-field up to the next multiple of its basic type. The bit-field is C's way of making bit-packing more ergonomic, a goal we accomplish with layout specifications. In our formalisation, each individual field of the bit-field does not stand alone as its own type; rather, we capture them all as a list of field-range pairs under one type, $\mathsf{bitfield}\left(\tau, \overline{f_i : r_i}\right)$.

```
struct {
    int a  : 8; int : 2; int b  : 12; int : 0;
    int c  : 4; int : 0;
    int d  : 20;
} str;
```

$$\mathsf{bitfield}(\ a : (\!|\, 0b, 8b \,|\!)\, ,\ b : (\!|\, 10b, 12b \,|\!)\, ,\ c : (\!|\, 1B, 4b \,|\!)\, ,\ d : (\!|\, 8B, 20b \,|\!)\ )$$
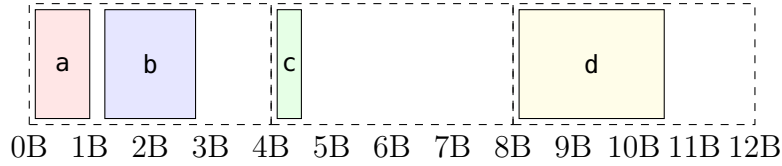


Figure 5: A bit-field, specified in C, shown in memory, and captured by our formalisation.

## Reducing to CR

Even this stripped-back formalisation includes more information than we need for reasoning about memory. We introduce a reduced type system CR (Figure 6), for two reasons: first, to minimize the complexity of our matching relation later, and second, to give the reader an intuitive sense of what should matter to us about these C types, and what can be discarded. In particular, note that CR features a set of basic memory constants $\rho$ which are not arbitrarily sized.

Now we establish a reduction $\rightsquigarrow$ from C to CR (Figure 7). Some reductions are trivial, but others provide insight. For example, VOIDRED codifies that, in accordance with the language specification, void types are layed out using the char memory representation. Many complex types drop information; for example, it doesn't matter what type a $\mathsf{pointer}\,(\tau)$ is pointing to, since all pointers use the same constant representation pointer. The arguments and return type of a $\mathsf{fun}\,(\overline{\tau_i}, \overline{\tau_j}, \tau_r)$ are also not important to us, as all functions are simply represented as pointers. It is important to remember that we are not checking the correctness of these types, which is what allows us to discard so much.

$$\begin{array}{lll}
\text{Basic memory types } \rho & ::= \text{pointer} \mid \text{char} & \text{(memory constants)} \\
& \mid \text{bool} \mid \text{short} \mid \text{int} \mid \text{long} \mid \text{long long} \\
& \mid \text{float} \mid \text{double} \mid \text{long double} \\
\text{Types} \qquad \tau' & ::= \rho & \text{(basic type } \rho \text{ aligned as } \tau') \\
& \mid \tau'[n] \\
& \mid \left\{ \overline{f_i : \tau'_i} \right\} \\
& \mid \left\langle \overline{f_i : \tau'_i} \right\rangle \\
& \mid \left| \tau', \overline{f_i : r_i} \right| \\
& \mid \tau'_1 \text{ aligned like } \tau'_2 \\
& \mid \text{as } \ell'
\end{array}$$

Figure 6: Reduced type system CR, containing only the information necessary for matching against memory layouts.

## Layouts in L

We now introduce a syntax for memory layouts in C (Figure 8). We'll refer to this layout language as L. It has slightly less syntax than DARGENT, as unions in C are not tagged, allowing them to be handled as simple memory sizes. One difference between L and DARGENT is the inclusion of the basic types $\rho$ as layouts, useful due to C's large number of primitive memory types. Also, in addition to placing one field after another, L includes syntax to place one field *before* another. Aside from these minor differences, its language design is very similar to DARGENT.

### Reducing to LR

Similar to the approach taken by Chen ([3]), we establish a core layout language LR on which we apply our semantics. Layouts in LR fundamentally consist of bit ranges describing the size and position of each piece of data, as well as structs containing a number of named fields. Endianness is still present in LR, with the endianness of the target machine represented by ME.

Reduction from L to LR (Figure 10) takes place in a context where:

- $Q$ is the environment of layout variables $\overline{l : \ell'}$.

- $L$ contains the layouts of the other fields present in the struct $\overline{f : \ell}$, used for relative positioning.

These rules simply desugar the syntax of L, reducing it to bit ranges. In ATRED, the absolute position is translated to an offset which is applied to the layout. In BEFORERED and AFTERRED, the name of the field $f$ is looked up in $L$, and the offset $o$ is calculated to position the layout either directly before or directly after the other field.

$$\boxed{B, Q \vdash \tau \leadsto \tau'}$$

$$\frac{B, Q \vdash \tau \leadsto \tau' \qquad B, Q \vdash \tau_a \leadsto \tau_a'}{B, Q \vdash \tau \,(\text{aligned like } \tau_a) \leadsto \tau' \text{ aligned like } \tau_a'}\textsc{AlignAttrRed}$$

$$\frac{B, Q \vdash \tau \leadsto \tau' \qquad \alpha \neq \text{aligned like } \tau_a}{B, Q \vdash \tau \,(\alpha) \leadsto \tau'}\textsc{DropAttrRed}$$

$$\frac{}{B, Q \vdash \text{void} \leadsto \text{char}}\textsc{VoidRed} \qquad \frac{}{B, Q \vdash \text{char}(\sigma) \leadsto \text{char}}\textsc{CharRed}$$

$$\frac{}{B, Q \vdash \text{int}(\sigma, \text{bool}) \leadsto \text{bool}}\textsc{BoolRed} \qquad \frac{}{B, Q \vdash \text{int}(\sigma, \text{short}) \leadsto \text{short}}\textsc{ShortRed}$$

$$\frac{}{B, Q \vdash \text{int}(\sigma, \text{default}) \leadsto \text{int}}\textsc{IntRed} \qquad \frac{}{B, Q \vdash \text{int}(\sigma, \text{long}) \leadsto \text{long}}\textsc{LongRed}$$

$$\frac{}{B, Q \vdash \text{int}(\sigma, \text{long long}) \leadsto \text{long long}}\textsc{LongLongRed}$$

$$\frac{}{B, Q \vdash \text{float} \leadsto \text{float}}\textsc{FloatRed}$$

$$\frac{}{B, Q \vdash \text{double}(\text{default}) \leadsto \text{double}}\textsc{DoubleRed}$$

$$\frac{}{B, Q \vdash \text{double}(\text{long}) \leadsto \text{long double}}\textsc{LongDoubleRed}$$

$$\frac{\tau \leadsto \tau'}{B, Q \vdash \text{array}\,(\tau, n) \leadsto \tau'\,[n]}\textsc{ArrayRed}$$

$$\frac{\text{for each } i : \tau_i \leadsto \tau_i'}{B, Q \vdash \text{struct}\,\left(\overline{f_i : \tau_i}\right) \leadsto \left\{\overline{f_i : \tau_i'}\right\}}\textsc{StructRed} \qquad \frac{\text{for each } i : \tau_i \leadsto \tau_i'}{B, Q \vdash \text{union}\,\left(\overline{f_i : \tau_i}\right) \leadsto \left\langle \overline{f_i : \tau_i'}\right\rangle}\textsc{UnionRed}$$

$$\frac{}{B, Q \vdash \text{pointer}\,(\tau) \leadsto \text{pointer}}\textsc{PointerRed} \qquad \frac{}{B, Q \vdash \text{enum}\,\left(\overline{f_i}\right) \leadsto \text{int}}\textsc{EnumRed}$$

$$\frac{B.\beta \leadsto \tau'}{B, Q \vdash \beta \leadsto \tau'}\textsc{TyVarRed} \qquad \frac{}{B, Q \vdash \text{fun}\,\left(\overline{\tau_i}, \overline{\tau_j}, \tau_r\right) \leadsto \text{pointer}}\textsc{FunRed}$$

$$\frac{}{B, Q \vdash \text{bitfield}\,\left(\tau, \overline{f_i : r_i}\right) \leadsto \left|\tau', \overline{f_i : r_i}\right|}\textsc{BitFieldRed}$$

$$\frac{}{B, Q \vdash \tau \text{ as } l \leadsto \text{as } Q.l}\textsc{AsRed}$$

Figure 7: Rules for reducing C to CR.

$$
\begin{array}{llr}
\text{Layouts} & \ell ::= s & \text{(memory blocks)} \\
& \mid\ \rho & \text{(memory constants)} \\
& \mid\ \ell \text{ at } s & \text{(absolute offsets)} \\
& \mid\ \ell \text{ after } f \mid \ell \text{ before } f & \text{(relative offsets)} \\
& \mid\ \ell \text{ using } \omega & \text{(endianness)} \\
& \mid\ l & \text{(variables)} \\
& \mid\ \left\{ \overline{f_i : \ell_i} \right\} & \text{(structs)} \\
\text{Endianness}\, \omega ::= \text{BE} \mid \text{LE}
\end{array}
$$

Figure 8: The syntax of our C layout language, L.

$$
\begin{array}{llll}
\text{Layouts} & \ell' & ::= & r_{\omega'} & \text{(primitive layouts)} \\
& & \mid & \left\{ \overline{f_i : \ell'_i} \right\} & \text{(structs)} \\
\text{Endianness} & \omega' & ::= & \text{BE} \mid \text{LE} \mid \text{ME} \\
\text{Field names} & & \ni & f
\end{array}
$$

$$\text{start, end, len} : \text{Layout} \to \text{Sizes}$$

$$\text{start}\left( \langle\!\langle\, o, s\, \rangle\!\rangle_{\omega'} \right) = o \qquad\qquad \text{end}\left( \langle\!\langle\, o, s\, \rangle\!\rangle_{\omega'} \right) = o + s$$

$$\text{start}\left( \left\{ \overline{f_i : \ell'_i} \right\} \right) = \min\left( \bigcup_i \text{start}\left( \ell'_i \right) \right) \qquad \text{end}\left( \left\{ \overline{f_i : \ell'_i} \right\} \right) = \max\left( \bigcup_i \text{end}\left( \ell'_i \right) \right)$$

$$\text{len}\left( \ell' \right) = \text{end}\left( \ell' \right) - \text{start}\left( \ell' \right)$$

$$\text{offset} : \text{Layouts} \to \text{Layouts}$$

$$\text{offset}\left( o, \langle\!\langle\, o', s\, \rangle\!\rangle_{\omega'} \right) = \langle\!\langle\, o + o', s\, \rangle\!\rangle_{\omega'}$$
$$\text{offset}\left( o, \left\{ \overline{f_i : \ell'_i} \right\} \right) = \left\{ \overline{f_i : \text{offset}\left( o, \ell'_i \right)} \right\}$$

Figure 9: Our core layout language, LR, and some basic operations.

## Static Semantics

We now consider the semantics of LR, starting with well-formedness rules (Figure 11). Here, we think of each layout as denoting a particular *allocation*, a set of bits which it occupies. If a layout considers the same bit to be a part of multiple fields, it is ill-formed. If an allocation can be found for a layout, it is well-formed; that is, $l \ \mathbf{wf} \stackrel{\text{def}}{=} \exists a.\ell' \asymp a$. Like Chen ([3]), we leave out the scope checking of layout variables.

Importantly, a layout's well-formedness does not guarantee that it will be able to match with any particular type. For this, we establish the *matching relation*. Before we can do this, we need to know the size of C's types, which is not a straightforward calculation.

To understand C sizes and alignment requirements, we formalise C's default scheme for laying out a particular type, which we call the *auto-matching relation*. The auto-matching

10

$$\boxed{Q, L \vdash \ell \rightsquigarrow \ell'}$$

$$\frac{0\mathsf{b} \vdash \rho \mathrel{\bar{\sim}} \ell'}{Q, L \vdash \rho \rightsquigarrow \ell'}\text{BASICTYPERED} \qquad \frac{}{Q, L \vdash s \rightsquigarrow \left(\!\left[\, 0\mathsf{b}, s \,\right]\!\right)_{\mathsf{ME}}}\text{SIZERED}$$

$$\frac{Q, L \vdash \ell \rightsquigarrow \ell'}{Q, L \vdash \ell \text{ at } s \rightsquigarrow \text{offset}\,(s, \ell')}\text{ATRED} \qquad \frac{}{Q, L \vdash l \rightsquigarrow Q.l}\text{VARRED}$$

$$\frac{\begin{array}{c} Q, L \vdash L.f \rightsquigarrow \ell'_f \\ Q, L \vdash \ell \rightsquigarrow \ell' \qquad o = \left[\text{end}\left(\ell'_f\right) + 1\right]\mathsf{b} \end{array}}{Q, L \vdash \ell \text{ after } f \rightsquigarrow \text{offset}\,(o, \ell')}\text{AFTERRED}$$

$$\frac{\begin{array}{c} Q, L \vdash L.f \rightsquigarrow \ell'_f \\ Q, L \vdash \ell \rightsquigarrow \ell' \qquad o = \left[\text{start}\left(\ell'_f\right) - \text{len}\left(\ell'\right)\right]\mathsf{b} \end{array}}{Q, L \vdash \ell \text{ before } f \rightsquigarrow \text{offset}\,(o, \ell')}\text{BEFORERED}$$

$$\frac{\text{for each } i\colon\; Q, \overline{f_i : \ell_i} \vdash \ell_i \rightsquigarrow \ell'_i}{Q, L \vdash \left\{\overline{f_i : \ell_i}\right\} \rightsquigarrow \left\{\overline{f_i : \ell'_i}\right\}}\text{STRUCTRED}$$

Figure 10: Rules for reducing layouts to core layouts.

11

$$\boxed{\ell' \asymp a}$$

$$\frac{}{(\!|\,o, s\,|\!)_{\omega'} \asymp \{\, \mathrm{bits}\,(o)\,, \mathrm{bits}\,(o) + 1, \cdots, \mathrm{bits}\,(o + s)\,\}} \textsc{RangeWF}$$

$$\frac{\text{for each } i\colon \ell'_i \asymp a_i \qquad \overline{a_i} \text{ disjoint}}{\{f_i : \ell'_i\} \asymp \bigcup_i a_i} \textsc{StructWF}$$

Figure 11: Well-formedness on core layouts.

relation $\ell' \mathbin{\bar\sim} \tau'$, means that, without our intervention, C's compiler will lay out a type $\tau'$ as $\ell'$ (Figure 12). These automatic layouts have some implementation-specific behaviour; in particular, the sizes and alignment restrictions on basic memory types $\rho$ vary between compilers, and often between target machines. Implementation-dependent behaviours have been  highlighted  in the figure. Because certain types will only align to multiples of certain sizes, the auto-matching relation is parameterised over $o$, the current offset of the type within the struct.

The fields of a struct are automatically layed out one after the other, and the whole struct is aligned according to the requirements of their strictest-aligned field (a complex alignment requirement represented in the rule by $\max (\cup_j o_j) \vdash \ell'_1 \mathbin{\bar\sim} \tau'_1$). Unions similarly use the strictest alignment among their fields, and they allocate as much space as their largest field requires. We simply treat them as a range,

Now we can establish the matching relation $\ell' \sim \tau'$, stating more generally that a layout $\ell'$ *can be* applied to a type $\tau'$ (Figure 13). According to RANGEMATCH, the alignment of types is overruled by our custom layout language, allowing any type to be placed anywhere by the user, provided that the layout itself is well-formed. In lieu of the tags present in COGENT, unions are simply matched with RANGEMATCH. Similarly, RANGEMATCH captures aligned types, and types that have been assigned a layout. Structs match when each field's type corresponds with a matching layout field. For bit-fields, the size of each field must be a match, but the fields can be positioned in any well-formed way.

$$\boxed{o \vdash \ell' \mathbin{\bar{\sim}} \tau'}$$

$$\frac{}{o \vdash (\!| \text{ aligned } (o, \text{ 8B }), \text{ 8B } |\!)_{\mathsf{ME}} \mathbin{\bar{\sim}} \mathsf{pointer}} \text{\textsc{PointerAuto}}$$

$$\frac{}{o \vdash (\!| \text{ aligned } (o, \text{ 1B }), \text{ 1B } |\!)_{\mathsf{ME}} \mathbin{\bar{\sim}} \mathsf{char}} \text{\textsc{CharAuto}}$$

$$\frac{}{o \vdash (\!| \text{ aligned } (o, \text{ 2B }), \text{ 2B } |\!)_{\mathsf{ME}} \mathbin{\bar{\sim}} \mathsf{short}} \text{\textsc{ShortAuto}}$$

$$\frac{}{o \vdash (\!| \text{ aligned } (o, \text{ 4B }), \text{ 4B } |\!)_{\mathsf{ME}} \mathbin{\bar{\sim}} \mathsf{int}} \text{\textsc{IntAuto}}$$

$$\frac{}{o \vdash (\!| \text{ aligned } (o, \text{ 4B }), \text{ 4B } |\!)_{\mathsf{ME}} \mathbin{\bar{\sim}} \mathsf{long}} \text{\textsc{LongAuto}}$$

$$\frac{}{o \vdash (\!| \text{ aligned } (o, \text{ 8B }), \text{ 8B } |\!)_{\mathsf{ME}} \mathbin{\bar{\sim}} \mathsf{long\ long}} \text{\textsc{LongLongAuto}}$$

$$\frac{}{o \vdash (\!| \text{ aligned } (o, \text{ 8B }), \text{ 8B } |\!)_{\mathsf{ME}} \mathbin{\bar{\sim}} \mathsf{float}} \text{\textsc{FloatAuto}}$$

$$\frac{}{o \vdash (\!| \text{ aligned } (o, \text{ 8B }), \text{ 8B } |\!)_{\mathsf{ME}} \mathbin{\bar{\sim}} \mathsf{double}} \text{\textsc{DoubleAuto}}$$

$$\frac{}{o \vdash (\!| \text{ aligned } (o, \text{ 8B }), \text{ 12B } |\!)_{\mathsf{ME}} \mathbin{\bar{\sim}} \mathsf{long\ double}} \text{\textsc{LongDoubleAuto}}$$

$$\frac{o \vdash \ell' \mathbin{\bar{\sim}} \tau'}{o \vdash (\!| \text{ start } (\ell'), \text{len} (\ell') \cdot n |\!)_{\mathsf{ME}} \mathbin{\bar{\sim}} \tau' [n]} \text{\textsc{ArrayAuto}} \qquad \frac{}{o \vdash \ell' \mathbin{\bar{\sim}} \mathsf{as}\ \ell'} \text{\textsc{AsAuto}}$$

$$\frac{\begin{array}{c} o \vdash \ell'_1 \mathbin{\bar{\sim}} \tau'_1 \qquad o \vdash \ell'_2 \mathbin{\bar{\sim}} \tau'_2 \\ p = \text{start} (\ell'_2) - \text{start} (\ell'_1) \end{array}}{o \vdash \text{offset} (p, \ell'_2) \mathbin{\bar{\sim}} \tau'_1 \ \mathsf{aligned\ like}\ \tau'_2} \text{\textsc{AlignedAuto}}$$

$$\frac{\begin{array}{c} \{\text{for each } j\colon o \vdash \ell'_j \mathbin{\bar{\sim}} \tau'_j \qquad \text{let } o_j = \text{start} (\ell'_j)\} \\ \max\left(\bigcup_j o_j\right) \vdash \ell'_1 \mathbin{\bar{\sim}} \tau'_1 \qquad \text{for all } i > 1\colon \text{end} (\ell'_{i-1}) \vdash \tau'_i \mathbin{\bar{\sim}} \ell'_i \end{array}}{o \vdash \overline{\{f_i : \ell'_i\}} \mathbin{\bar{\sim}} \overline{\{f_i : \tau'_i\}}} \text{\textsc{StructAuto}}$$

$$\frac{\text{for all } i\colon o \vdash \ell'_i \mathbin{\bar{\sim}} \tau'_i}{o \vdash (\!| \max\left(\bigcup_i \text{start} (\overline{\ell'_i})\right), \max\left(\bigcup_i \text{len} (\overline{\ell'_i})\right) |\!)_{\mathsf{ME}} \mathbin{\bar{\sim}} \left\langle \overline{f_i : \tau'_i} \right\rangle} \text{\textsc{UnionAuto}}$$

$$\frac{}{\overline{\{f_i : (r_i)_{\mathsf{ME}}\}} \mathbin{\bar{\sim}} \left| \tau', \overline{f_i : r_i} \right|} \text{\textsc{BitfieldAuto}}$$

Figure 12: Auto-matching relation $\bar{\sim}$ between CR and LR.

$$\boxed{\ell' \sim \tau'}$$

$$\frac{\ell' \ \bar{\sim} \ \tau'}{(\!| \, o, \mathrm{len}\,(\ell') \, |\!)_{\omega'} \sim \tau'}\textsc{RangeMatch}$$

$$\frac{\text{for each } i\colon \ell'_i \sim \tau'_i}{\overline{\left\{ f_i : \ell'_i \right\}} \sim \overline{\left\{ f_i : \tau'_i \right\}}}\textsc{StructMatch}$$

$$\frac{}{\overline{\left\{ f_i : (\!| \, o_i, \mathrm{len}\,(r_i) \, |\!)_{\mathsf{ME}} \right\}} \sim \left| \tau', \overline{f_i : r_i} \right|}\textsc{BitfieldMatch}$$

Figure 13: Matching relation between reduced layouts and CR types.

## Layout Variable Declaration

We have not addressed the declaration of layout variables; that is, the addition of new $l : \ell'$ entries into our layout variable environment $Q$. For this, we insert layout variable declarations into C's statement syntax. We don't give a full account of the C statement syntax since it comprises a full chapter in the specification ([6]), but we simply append the syntax as follows:

$$\text{Statements } s ::= \cdots \mid \ l = \ell;$$

Next we give a denotational semantics of these layout variable declarations. Layout variable declarations are thought of as an insertion into the environment $Q$, where the reduction to a core layout $\ell'$ happens at the point of declaration, taking in the layout variable context $Q$, and the empty set for the field context $L$ (since it is populated within the reduction, by the StructRed rule).

$$Q \ [\![ \, l = \ell; \, ]\!] = l : \ell'$$
$$\text{where } Q, \{\} \vdash \ell \rightsquigarrow \ell'$$

## Implementation in C

The focus of this work is purely pen-and-paper, exploring the feasibility of layout specifications interacting with C's type system. As such, we do not thoroughly explore implementation options. DARGENT, as an extension of COGENT, is implemented as a transpilation to C, with custom layouts in DARGENT transpiling to singleton C structs, with generated getter and setter code ([1]). A similar approach could be taken here, with an extended language C+L transpiling down to C, generating low-level bit-twiddling code to get and set values with custom layouts. A more direct approach would be to include layout specification as a part of the C language, with compilers handling the layouts in their compilation. The language design implications of this would be significant.

# 5    Discussion

The following observations arose from studying the C and HASKELL language specifications in preparation for writing this report, research which was more difficult than expected. In the case of C, an understanding of type layout comes partially from the C standard, and partially from GCC documentation. For HASKELL, the language specification leaves memory layout very open - GHC documentation is the primary resource, although conclusive information about even the basic memory layout strategies employed by GHC is buried deep within compiler wikis. The spuriousness and ambiguity of the current state of memory layout formalisation is the context for this extended line of inquiry.

## Limitations of C Types

Although our layout language for C provides a good basis for future work in this direction, its limitations are clear. Our attempts to create freedom of type representation ran into limitations imposed by the standard. The C standard sits halfway between enforcing restrictions on memory layout so that programmers can write bit-manipulating code, and leaving the details open to implementation so that compilers can optimise. This former part is problematic for us; we want programmers to have control over the bit representation not by simply formalising one fixed representation, but rather through flexibly-represented types which they can specify themselves - made impossible when the language design restricts memory layout. The latter part is also problematic, since leaving the details of memory layout completely unwritten means we rely on compiler documentation, which is often informal or incomplete. Then, any code we write will depend on the behaviour of some particular compiler or machine - behaviour which could change under our feet.

One immediate solution would be to add new types to C whose length is user-specified, like the U$n$ family of types in COGENT. That way, fields could be created in structs which match custom-size layouts. However, with this solution not only do the type signature and the layout specification contain duplicate information, but it also progresses us further down the path of inseperability of types from their memory representations. To approach the problem differently, we must consider the semantics of types.

## Disentangling Language and Memory

Ideally, an object's type conveys semantic information - an integer or a boolean is called as such not because of the way it looks in memory, but because of the *range of things it represents*. Thus, a language should be formalisable in layout-agnostic terms, so that it is truly portable and maximally flexible. Reynolds gives an account of types discussing their semantics in these terms, with type semantics like $[\![\ \mathsf{Int}\ ]\!] = \mathbb{Z}_\perp$ ([7]).

The question then arises: how can programmers write code which takes advantage of facts about memory representation, if it isn't given by the language specification? Here, we need new specification. In particular, it is necessary for language implementers to formalise their treatment of memory for any given type, and for this implementation specification to modularly extend or refine the language semantics, allowing for different compilers to be

properly compared. What follows from this is a disentanglement of language semantics from memory.

We will go through one toy example to communicate the idea. Suppose a language A has a type SmallNat. The language specification gives a domain for this type as follows:

$$[\![ \, \text{SmallNat} \, ]\!] = [0, 255]$$

Note the domain of SmallInt is *not* the entire domain of $\mathbb{N}$, and that this design decision is informed by practicality of implementation - however, it still stops short of laying out memory. What is important is that the domain is semantic, referring to mathematical objects outside of the language.

The implementation specification then describes read and write functions for the type - that is, functions from the domain of the type in A, to a domain made up of memory units which are, in this case, binary.

$$\text{Unit} = \mathbb{B}$$

$$\text{read} : \mathbb{B}^8 \to [0, 255] \qquad\qquad \text{write} : [0, 255] \to \mathbb{B}^8$$

$$\text{read}\,(b) = \sum_{i=0}^{7} b_i \cdot 2^i \qquad\qquad \text{write}\,(n) = \text{write}\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + \langle n \bmod 2 \rangle$$

We could instead specify read and write functions for a most-significant bit ordering, an equivalently compact representation on the same hardware but with different read-write semantics:

$$\text{read} : \mathbb{B}^8 \to [0, 255] \qquad\qquad \text{write} : [0, 255] \to \mathbb{B}^8$$

$$\text{read}\,(b) = \sum_{i=0}^{7} b_i \cdot 2^{7-i} \qquad\qquad \text{write}\,(n) = \langle n \bmod 2 \rangle + \text{write}\left(\left\lfloor \frac{n}{2} \right\rfloor\right)$$

Suppose, for some reason, we want to compile A on a machine which stores its memory in ternary units, $\mathbb{T} = \{\, 0, 1, 2 \,\}$. We can still establish functions for this without changing the language design, using six trits instead of eight bits:

$$\text{Unit} = \mathbb{T}$$

$$\text{read} : \mathbb{T}^6 \rightharpoonup [0, 255] \qquad\qquad \text{write} : [0, 255] \to \mathbb{T}^6$$

$$\text{read}\,(b) = \sum_{i=0}^{5} b_i \cdot 3^{5-i} \qquad\qquad \text{write}\,(n) = \langle n \bmod 3 \rangle + \text{write}\left(\left\lfloor \frac{n}{3} \right\rfloor\right)$$

The point of this example is how the same language specification can be refined to suit several different implementation contexts without adjustment, because it is defined without reference

to memory. Thus, anything proven about A holds across these different environments, with no danger of undefined behaviour.

From this foundation, we could equip the programmer with surface syntax exposing the read and write behaviour of types, allowing them to reason about abstract datatypes while retaining arbitrarily fine control over how those types are represented, going beyond just the limited degree of control we offer in our proposed language L. A programmer in special circumstances may need to represent their integers using an entirely different representation, such as a different bit ordering; with this proposed approach, they would simply need to write reliable read and write functions, and then the compiler would allow them to use their deeply-customised integers as though they were regular integers (albeit potentially with less efficient performance, depending on their representation).

If formally describing the memory representation of every type sounds too time-consuming, consider that once basic types as well as product and sum types have been specified, many other commonly-used types will immediately follow. Importantly, if an existing language is to be ported to new hardware, its language semantics do not need to change, in the way that C's semantics have needed to - famously, advancements in processor hardware motivated the introduction of the long type because it was the simplest way to avoid changing specification of the int type ([6]). Many of C's complexities may have been avoided had the designers separated layout from their core language and treated it as an implementation detail, to be formalised elsewhere in a modular fashion. A universal, machine-readable language of memory layout specification could be established, allowing language implementations to draw on common libraries of memory representations. This research direction has a lot of promise; further investigation should be done into comparative study of memory layouts in various languages, to establish a common language of memory representation and finally disentangle types from layout.

# 6 Conclusion

Our investigations into data layout languages have resulted in a successful extension to DARGENT's formalisation, accounting for recursive types and treating them as pointers. In addition, we have developed an entirely new layout syntax for C, allowing for the manipulation of struct fields, enabling programmers to easily rearrange and pack their data. Aside from following the steps established by Chen in layout language design ([3]), we also outline C's automatic approach to layout, which provides an understanding of the its type system specificities, particularly with regards to padding and sizing. Finally, we outline possible research directions in the separation of a memory layouts from language semantics, as a new project pursuing portability and reliability.

## Acknowledgements

# References

[1] Chen, Z., Lafont, A., O'Connor, L., Keller, G., McLaughlin, C., Jackson, V., and Rizkallah, C. (jan, 2023) Dargent: A Silver Bullet for Verified Data Layout Refinement. *Proc. ACM Program. Lang.,* **7**(POPL).

[2] Baudon, T., Radanne, G., and Gonnord, L. (aug, 2023) Bit-Stealing Made Legal: Compilation for Custom Memory Representations of Algebraic Data Types. *Proc. ACM Program. Lang.,* **7**(ICFP).

[3] Chen, Z. Towards A Practical High-Assurance Systems Programming Language PhD thesis UNSW Sydney (2023).

[4] Murray, E. Recursive types for cogent PhD thesis Honours thesis, UNSW Sydney, Computer Science & Engineering (2019).

[5] Leroy, X. (2009) Formal verification of a realistic compiler. *Communications of the ACM,* **52**(7), 107–115.

[6] ISO (June, 2018) ISO/IEC 9899:2018 Information technology — Programming languages — C, , fourth edition.

[7] Reynolds, J. C. What do types mean? — From intrinsic to extrinsic semantics pp. 309–327 Springer New York New York, NY (2003).