

Bit Level Types: Syntax and Semantics

David T Eger (eger@cs.cmu.edu)
Carnegie Mellon University
School of Computer Science
Pittsburgh, PA 15213

September 22, 2005

0 To Do

Herein is a draft of the semantics for BLT. This draft is missing some key features from full-blown BLT, notably the constructs `implicit`, `dontcare`, `asbits`, `wildcards`, the locatives `@.!` and `@!`, support for opaque types such as IEEE 754 floats, and the specific rules for type embedding backends for (C/SMLNJ/etc.).

1 Overview

Suppose we are given a document which describes the Shockwave Flash File Format (SWF). This document will go to great lengths describing how to interpret the sequence of bits in a shockwave file as structured data. Certain sequences of bits will translate to strings of text. Certain sequences of bits will translate to data structures representing rectangles. Certain sequences of bits are raw data that must be run through a decompression algorithm before their true structure can be interpreted.

Suppose we would like to write a program in C which manipulates SWF files. We must take the prose and figures which are the SWF specification, and write down an interpretation of them as structured data in C. Some portions of the specification may not map cleanly to C. For instance, a rectangle is described in Flash as a structure with five fields:

```
RECT =  
struct  
    n : uint(5),  
    xmin : uint(n),  
    xmax : uint(n),  
    ymin : uint(n),  
    ymax : uint(n)  
tcurts
```

BLT spec of an SWF RECT

```
typedef struct {  
    unsigned char n;  
    unsigned int xmin;  
    unsigned int xmax;  
    unsigned int ymin;  
    unsigned int ymax;  
} RECT;
```

A Possible Interpretation in C

The first (5-bit) integer field specifies how many bits each of the following four integer fields will take in the bit stream. From this we can infer that none of the fields in this structure will consume more than 32 bits, and so on a 32 bit platform it is safe to map SWF's notion of a rectangle to the C struct in the figure on the right. However, our interpretation of a SWF RECT in C *is not* a real SWF RECT. The *definition* of what a SWF RECT is is defined exactly by the SWF specification, which defines the *file format*.

Thus, there is often a *loss of information* about the structure of data when we interpret said data in a target language such as C. According to the SWF specification, the field *n* is only allowed to have values between 0 and 31. Our embedding in C, however, allows *n* to have any value between `'\x00'` and `'\0xFF'`. In our embedding

in C, we have also forgotten the dependency between the field n and the fields $xmin$, $xmax$, $ymin$, and $ymax$. In a value of type SWF RECT, the latter four fields must each be an integer in the range 0 and $2^n - 1$. In our embedding in C, we know only that each must be between 0 and $2^{32} - 1$.

Nonetheless, we can create an approximation of what a SWF RECT is within C in terms of C's data types. We must then write marshalling code to convert between the data on disk and our approximation of it in C.

As we parse, we must be careful to interpret the data correctly: accounting for bit field issues, dependency issues, and validation conditions on the data. When we want to write the data back out to disk, we need to ensure that the invariants on a Flash File which are not captured by the C type system have been preserved.

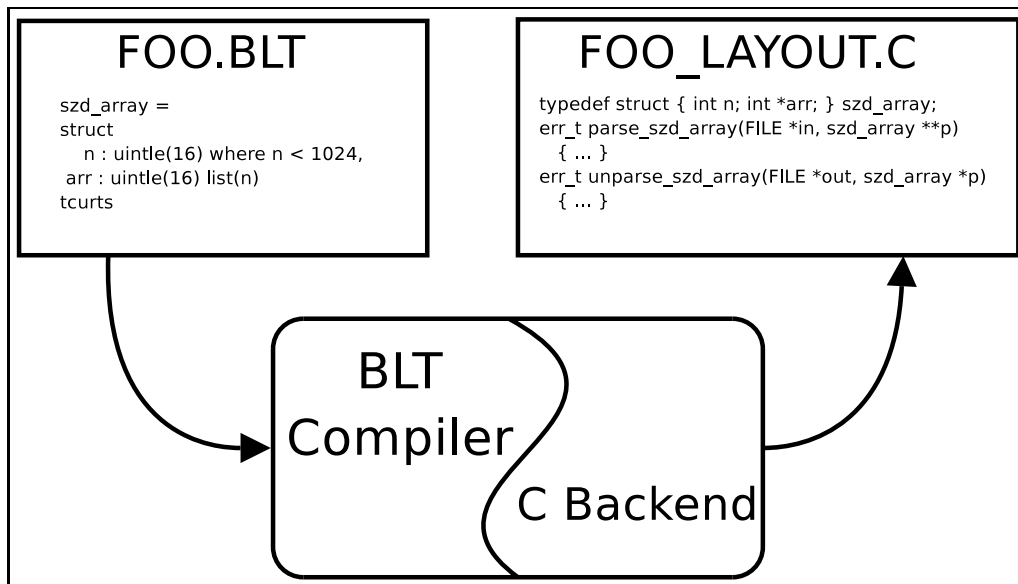
At the end of the day, we will have written two C functions:

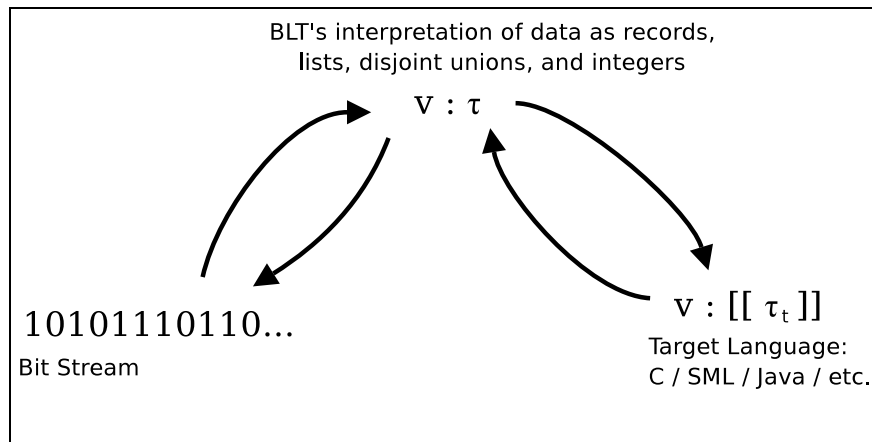
```
err_t parse_RECT(FILE * f, RECT **r);
err_t unparse_RECT(FILE * f, RECT *r);
```

The goal of BLT is to take the burden of interpreting low-level data off the shoulders of an application writer. Instead of laboriously translating a two-hundred page file format specification into parsing code, we provide a *specification language* (BLT). Each BLT specification specifies how a certain type of structured data is laid out on disk. The types of data we can specify are integers, records, disjoint unions, and lists. We provide a small set of terms for integer layout types and combinators for constructing layout terms for structured data out of smaller layout terms. A BLT compiler has a backend for a *target language* (be it C, SML, Java, etc.) which translates an abstract type in BLT to an appropriate concrete interpretation in that target language. For instance, the mapping of BLT types to C looks as follows:

BLT type	int	$\{v_i : \tau_i\}_{i=1}^n$	$[v_i : \tau_i]_{i=1}^n$	τ list
<i>C type</i>	One of : (un)signed char (un)signed short (un)signed int (un)signed long BIGNUM	<code>struct{([τ_i] v_i);}_{i=1}^n</code> ;	<code>struct{ int dm; union{([τ_i] v_i);}_{i=1}^n}; };</code>	<code>struct{ int len; [[τ]] *arr; };</code>

A BLT compiler will translate a data layout specification into appropriate parse and unparse functions in the programmer's target language.





An important design decision to mention before we go further is that our language is based on the presumption that the data we wish to parse can be parsed *predictively*. That is, we are working with some variant on the idea of LL(k) grammars: once we've looked at the next few bits of the incoming bit stream and decided on a branch of a union or whether a greedy list terminates, there's no need to backtrack.

We want this behavior so that we can parse streams in pieces. As an example, let's suppose we are writing a specification for the data sent over a socket in an X Protocol Session. We may describe this data as an XHandshake packet followed by a sequence of XRequests followed optionally by an XShutdown packet. This correctly describes a session between an X Client and an X Server. However, if we're using BLT to make a monitoring or fault injection tool, parsing the entire session is no good to us. We want the ability to change individual X Requests as they are made, and the X Protocol is designed so we can do this: each request, handshake, and shutdown packet can be parsed without knowing about any data that follows it.

One of our goals in designing BLT, then, will be to be able to provide warnings to the specification writer when his specification contains an ambiguous parse: one where the predictive parse is not the only parse.

2 Well Formed BLT Terms and Types

BLT is a language for concisely describing the layout of a data structure as a string of bits. Those terms which describe a data type and its layout are of type ‘ α layout.’ When compiled, such terms will be exposed to the programmer in a form reminiscent of combinator parsers: $\{\text{parse} : \text{bit string} \rightarrow \llbracket \alpha \rrbracket * \text{bit string}, \text{unparse} : \llbracket \alpha \rrbracket \rightarrow \text{bit string}\}$, where $\llbracket \alpha \rrbracket$ is the translation of a BLT type into the target language’s native type system.

$\tau ::=$	int	Integer (Unbounded)
	$\{l_i : \tau_i\}_{i=1}^n$	Record
	$[l_i : \tau_i]_{i=1}^n$	Disjoint Sum
	τ list	List of elements of type τ
	τ layout	Layout of a τ on a bitstream
	$\tau_1 \rightarrow \tau_2$	Function
$e ::=$	$\lambda v : \tau_1. e_2$	Lambda Abstraction
	$\lambda v \lambda e_1. e_2$	Lambda Abstraction with Type – Directed Sugar
	v	Reference A Bound Variable
	k	Konstant (see below)
	$e_1(e_2)$	Function Application
	e_i where e_2	Type Refinement/Validation Condition
	e_1 asa e_2	Checked Coercion
	$\text{struct}(v_i : e_i)_{i=1}^n \text{tcurts}$	Dependent Record Layout (no padding)
	$\text{union}(v_i : e_i)_{i=1}^n \text{noinu}$	Simple Sum (e_i need ensure representation discriminability)
	$e_1 \text{list}(e_2)$	List of e_1 ’s of Length e_2 (no padding)
	$e_1 \text{list}$	Greedy e_1 List (no padding)
	$\text{case } e \text{ of } (v_i(v) \Rightarrow e_i)_{i=1}^n$	Elimination Form For Sums
	$e_1.v$	Elimination Form For Records

For convenience of presentation we will define $\text{bool} \stackrel{\text{def}}{=} \{\text{true} : \{\}, \text{false} : \{\}\}$ and the syntax “if e then e_1 else e_2 ” as sugar for “case e of $\text{true} \Rightarrow e_1, \text{false} \Rightarrow e_2$.” We’ll also represent tuples in the common way: a tuple is just a record with natural number labelled fields $1 \dots n$.

	<i>Type</i>	<i>Description</i>
$k ::=$	true false	Boolean Constant
	AND OR	Boolean Connectives
	NOT	Boolean Negation
	$\dots, -1, 0, 1, \dots$	Integer Constant
	length	$\forall \alpha. \alpha \text{ list} \rightarrow \text{int}$ Give the Length of a list
	nth	$\forall \alpha. \alpha \text{ list} * \text{int} \rightarrow \alpha$ Project the n^{th} Element of a list
	$= \neq < > \leq \geq $	Integer Comparison
	$+ * / \% -$	Integer Operations
	$\& \ll \gg$	“Bitwise” Operators (partial! see later)
	uint	$\text{int} \rightarrow \text{int layout}$ n bit unsigned integer
	uintle	$\text{int} \rightarrow \text{int layout}$ n bit unsigned little endian integer
	sint	$\text{int} \rightarrow \text{int layout}$ n bit two’s complement signed integer
	sintle	$\text{int} \rightarrow \text{int layout}$ n bit two’s complement signed little endian integer

For syntactic convenience, we shall provide infix notation for the common arithmetic operators.

One of the oddities to note is that we have a number of symbols that are or should arguably be colons (see table below). This is due to the nature of the type τ layout.

We can view terms of type τ layout from a couple of viewpoints. First, we can view a term of type τ layout as specifying a marshaller for objects of type τ which implements one of several possible layout strategies.

Secondly, we can view terms of type τ layout as specifying refinements of τ according to its dependencies and where-clauses – things not accounted for by the BLT type system. (These constraints we will reason about using a form of Abstract Interpretation.) Both of these properties of layout terms are needed in order to formalize the prose that comprises the original specifications for binary formats.

We therefore view the syntactic constructs `struct ... tcurts`, `union ... noinu`, `asa`, `list` and `where` as marshaller combinators, mirroring the structure of the type system.

Having explained the correspondence between marshaller combinators (terms) and types, let us now explain the meaning of each colon wherever it may occur.

Colon-like Symbol	Example Usage	Meaning
:	<code>{a : int, b : int}</code>	Type of a field of a Record Type
:	<code>[a : int, b : int]</code>	Type of a labeled variant of a Union Type
:	<code>term : typ</code>	The BLT term <code>term</code> has the type <code>typ</code>
:	<code>struct a : uint(8) tcurts</code>	Part of the syntax of the <code>struct</code> record layout introduction form: This field should be marshalled by the layout term on the right
:	<code>union a : uint(8) noinu</code>	Part of the syntax of the <code>union</code> sum layout introduction form: This branch should be marshalled by the layout term on the right
:	<code>$\lambda v : \tau_1. e_2$</code>	Lambda Abstraction
}	<code>$\lambda v \{ e_1. e_2$</code>	Type-Directed Sugar. <code>e_1</code> must be of type τ layout, and this is equivalent to <code>$\lambda v : \tau. e_2$</code>

Let's start with an example:

```
struct a : uint(8), b : uint(16) tcurts
```

is of type “`{a : int, b : int}` layout”. `struct` constructs itsmarshallers by sequencing a series ofmarshallers, passing the partially parsed results forward during the parse phase, and passing the complete parsed structure to each marshaller during the unparse phase. Therefore, the term to the right of “`a :`” which is “`uint(8)`” needs to be a layout type, which it is.

Most of our table should be familiar, but our syntactic sugar “squiggle” may need some explanation. As illustrated in the extended example in Section 5, BLT specifications consist of a sequence of declarations. It is common for a layout term (in the example, `cu16`) to depend on something previously parsed (in the example, an XEndianness field). Therefore, it is natural to want to indicate an argument by simply naming the layout term which parsed the data instead of writing the BLT type of the parsed data explicitly.

Nonetheless, we may occasionally want to write a plain lambda, which we illustrate by:

```
min =  $\lambda x : \text{int}. \lambda y : \text{int}. \text{if } x < y \text{ then } x \text{ else } y$ 
```

3 Type Equality

$$\frac{}{\text{int} =_{\tau} \text{int}} \text{EQT-INT} \quad \frac{m = n \quad (\tau_i =_{\tau} \sigma_i \text{ and } l_i = m_i \text{ (as strings)})_{i=1\dots n}}{\{l_i : \tau_i\}_{i=1}^m =_{\tau} \{m_i : \sigma_i\}_{i=1}^n} \text{EQT-RCD}$$

$$\frac{m = n \quad (\tau_i =_{\tau} \sigma_i \text{ and } l_i = m_i \text{ (as strings)})_{i=1\dots n}}{[l_i : \tau_i]_{i=1}^m =_{\tau} [m_i : \sigma_i]_{i=1}^n} \text{EQT-SUM}$$

$$\frac{\tau =_{\tau} \sigma}{\tau \text{ list} =_{\tau} \sigma \text{ list}} \text{EQT-LIST}$$

$$\frac{\tau =_{\tau} \sigma}{\tau \text{ layout} =_{\tau} \sigma \text{ layout}} \text{EQT-LAYOUT}$$

$$\frac{\tau_1 =_{\tau} \tau_3 \quad \tau_2 =_{\tau} \tau_4}{\tau_1 \rightarrow \tau_2 =_{\tau} \tau_3 \rightarrow \tau_4} \text{EQT-FN}$$

4 Static Semantics

For each of the constants k above, we have a rule saying that in any context, said term is of the type we specified in the table. The constants are not all truly constants per sé, as for instance our BLT type system does not account for universal quantifiers. We account for terms like `length` and `nth`, then, using the same sort of “bolted on” approach that Pierce uses in the early chapters of TAPL. The rules for our constants follow from their types, so we omit them for brevity.

$$\begin{array}{c}
\frac{\Gamma, \nu : \tau_1 \vdash e_2 : \tau_2}{\Gamma \vdash (\lambda \nu : \tau_1. e_2) : \tau_1 \rightarrow \tau_2} \text{S-ABST} \quad \frac{\Gamma \vdash e_1 : \tau_1 \text{ layout} \quad \Gamma, \nu : \tau_1 \vdash e_2 : \tau_2}{\Gamma \vdash (\lambda \nu \lambda e_1. e_2) : \tau_1 \rightarrow \tau_2} \text{S-ABSLAYOUT} \\
\\
\frac{\Gamma \vdash e_1 : \tau_1 \rightarrow \tau \quad \Gamma \vdash e_2 : \tau_1}{\Gamma \vdash e_1(e_2) : \tau} \text{S-APP} \\
\\
\frac{\Gamma \vdash e : [l_i : \tau_i]_{i=1}^n \quad (l_i = \nu_i \text{ and } \Gamma, \nu : \tau_i \vdash e_i : \tau)_{i=1}^n}{\Gamma \vdash \text{case } e \text{ of } (\nu_i(\nu) = > e_i)_{i=1}^n : \tau} \text{S-ELIMUNION} \\
\\
\frac{\Gamma \vdash \text{foo} : \{ \nu_i : \tau_i \}_{i=1}^n}{\Gamma \vdash \text{foo}.\nu_k : \tau_k} \text{S-ELIMRCD} \\
\\
\frac{\Gamma \vdash e_1 : \alpha \text{ layout} \quad \Gamma, \$! : \alpha \vdash e_2 : \text{bool}}{\Gamma \vdash e_1 \text{ where } e_2 : \alpha \text{ layout}} \text{S-REFINE} \quad \frac{\Gamma \vdash e_1 : \alpha \text{ layout} \quad \Gamma \vdash e_2 : \alpha_2 \text{ layout}}{\Gamma \vdash e_1 \text{ asa } e_2 : \alpha_2 \text{ layout}} \text{S-REINTERPRET} \\
\\
\frac{\Gamma \vdash e_1 : \alpha \text{ layout}}{\Gamma \vdash e_1 \text{ list} : \alpha \text{ list layout}} \text{S-LIST} \quad \frac{\Gamma \vdash e_1 : \alpha \text{ layout} \quad \Gamma \vdash e_2 : \text{int}}{\Gamma \vdash e_1 \text{ list}(e_2) : \alpha \text{ list layout}} \text{S-NLIST} \\
\\
\frac{\Gamma \vdash e_i : \tau_i \text{ layout} \quad \text{for } i \in [1 \dots n] \quad \nu_i \text{ all distinct}}{\Gamma \vdash \text{union}(\nu_i : e_i)_{i=1}^n \text{ noinu} : [\nu_i : \tau_i]_{i=1}^n \text{ layout}} \text{S-UNION} \\
\\
\frac{}{\text{struct } \text{tcurts} : \{ \}_{i=1}^0 \text{ layout}} \text{S-STRUCTBASE} \\
\\
\frac{\Gamma \vdash \text{struct } (\nu_i : e_i)_{i=1}^{n-1} \text{ tcurts} : \{ \nu_i : \tau_i \}_{i=1}^{n-1} \quad \Gamma, (\nu_i : \tau_i)_{i=1}^{n-1} \vdash e_n : \tau_n \text{ layout}}{\Gamma \vdash \text{struct } (\nu_i : e_i)_{i=1}^n \text{ tcurts} : \{ \nu_i : \tau_i \}_{i=1}^n \text{ layout}} \text{S-STRUCTSTEP}
\end{array}$$

Static Semantics for BLT (sans rules for Konstants)

The first two rules of our static semantics, S-ABST and S-ABSLAYOUT re-iterate the distinction between our two colons at the term level. S-APP is straight-forward, but let us walk through the next four rules one at a time, hinting at our operational semantics. S-REFINE would be used to typecheck a term such as:

`uint(8) where $! > 5`

This term is of type `int layout`, and is different from `uint(8)` in that the integer represented (that is, the integer resulting from parsing or provided to unparse) must be greater than 5. Thus where-clauses represent refinements or validation conditions on data (i.e. boolean predicates that must be satisfied). The symbol “\$!” should be read as a “this” symbol, bound in the body of e_2 . A derivation for the above example can be worked as follows (rules for Konstants elided):

$$\frac{\cdot \vdash \text{uint} : \text{int} \rightarrow \text{int layout} \quad \cdot \vdash 8 : \text{int}}{\cdot \vdash \text{uint}(8) : \text{int layout}} \text{S-APP} \quad \frac{\$! : \text{int} \vdash (\$,5) : \text{int} * \text{int} \quad \$! : \text{int} \vdash > : \text{int} * \text{int} \rightarrow \text{bool}}{\$! : \text{int} \vdash \$! > 5 : \text{bool}} \text{S-APP}$$

$$\frac{\cdot \vdash \text{uint}(8) : \text{int layout} \quad \$! : \text{int} \vdash \$! > 5 : \text{bool}}{\cdot \vdash \text{uint}(8) \text{ where } \$! > 5 : \text{int layout}} \text{S-REFINE}$$

A Re-Interpreted type is a form of type-cast. The use of this construct is two-fold: first to guarantee that the low-level bits may be interpreted as a value of the type specified by e_1 , and secondly, e_1 determines the size of the prefix of the bitstream being parsed that must be interpreted as an e_2 .

We have two layout constructors for lists. The first simply takes an α layout and uses it greedily in the parse stage to read in a list of α 's. The second reads a list of α 's of some integer length e_2 .

unions represent the layout of disjoint union types. This is one point where the fact that our system is based on *predictive* parsing shows especially. Parsers for union types try each τ_i layout in turn until one successfully parses the prefix to the bit stream, and *does not ever try a different branch later on the same prefix*. To ensure round-trip preservation of data types, the specification writer needs to determine that the branches of a union encode to different bit strings.

The rules for constructing dependent records are fairly straight-forward: later fields' marshallers can be determined using the values parsed as earlier fields.

5 An Extended Example

Let's finish off this section by illustrating the first fragment of the X Protocol.

u8	=	uint(8)
XEndianness	=	union BIG : u8 where \$! = 'B', LITTLE : u8 where \$! = 'l' noinu
cu16	=	$\lambda c \{ XEndianness.$ case c of BIG(v) => uint(16) LITTLE(v) => uintle(16)
padto	=	$\lambda n : \text{int} . \lambda \text{sz_align} : \text{int}.$ if $n \% \text{sz_align} = 0$ then 0 else $n + (\text{sz_align} - n \% \text{sz_align})$
XClientHandshake	=	struct byte_order aka c : XEndianness, unused : u8, protocol_major_version : cu16(c) where \$! = 11, protocol_minor_version : cu16(c), n : cu16(c), d : cu16(c), protocol_name : u8 list(n), p0 : u8 list(padto(n)(4)), protocol_data : u8 list(d), p1 : u8 list(padto(d)(4)) tcurts

The above specification¹ is elaborated to have the following types:

u8 : int layout

¹BLT Specifications consist of a list of named expressions which are to be compiled and provided to the application programmer as a library. The two features we use here which we have not explained heretofore are parentheses within arithmetic expressions for grouping, and the 'aka' construct, which provides an alias for an otherwise long field name within a structure declaration.

XEndianness : [BIG : int, LITTLE : int] layout

cu16 : [BIG : int, LITTLE : int] → int layout

padto : int → int → int

XClientHandshake : { byte_order : [BIG : int, LITTLE : int], unused : int,
protocol_major_version : int,
protocol_minor_version : int,
n : int, d : int,
protocol_name : int list, p0 : int list,
protocol_data : int list, p1 : int list } layout

Note that in the definition of XClientHandshake, even though “uint(16)” and “uintle(16)” are not the same, because they are both of type “int layout,” they can be used as alternate branches of a case. One could similarly have “uint(8)” and “uintle(32) where \$! ≠ 56” as branches of a case statement, but the sort of dependencies BLT can express are only of this sort (those for which the parsed data is the same BLT type). The following BLT specification would not type-check since the branches of the if statement are of different types: [BIG : int, LITTLE : int] layout versus int layout.

$\lambda n : \text{int. if } n > 0 \text{ then XEndianness else uint}(32)$

6 Operational Semantics

One of the things that may seem strange about our Static Semantics is that there are basically no elimination forms for layout terms. Sure, we can use layout terms to construct bigger layout terms, but there’s no direct way to use these terms within BLT. This is by design. Layout terms are only actually used in their compiled form from within the target language to marshal structured data.

Our presentation of BLT’s Operational Semantics is for definitional purposes, and as such have left optimizations to the implementor. Nonetheless, there are some fundamental issues that we will address later which we gloss over here: those of choosing appropriate representations for data. For instance, we would like to interpret a uint(8) list in C as a struct {int len; char * arr}; instead of as a struct {int len; BIGNUM * arr};. For now, we ignore this detail and operate in an idealized setting where all integers in our computations have the semantics of BIGNUMS. We provide our operational semantics by means of a canonical compilation of a BLT expression into an SML library. These operational semantics mirror what an actual implementation would look like. We avoid the detail of reifying an abstract data structure of a BLT type into the target language by operating on native SML data explicitly.

XXX See attached unhygenic macro-ish SML. Can I use Aleks’s formalism here?

7 Operational Semantics (take 2)

Here, we provide interpretation of a BLT layout terms t by parse and unparse relations they define on strings of bits s :

$$\begin{array}{ll} s \cdot s', e \xrightarrow{p} v, s' & e \text{ parses the prefix } s \text{ to the value } v \text{ leaving } s' \\ v, e \xrightarrow{u} s & e \text{ unparses the value } v \text{ to the string } s \\ e_1 \mapsto e_2 & e_1 \text{ steps to } e_2 \text{ without affecting the state} \\ e_1 \rightarrow e_2 & \text{Transitive closure of } \mapsto \end{array}$$

$$\frac{e_1 \mapsto e'_1}{e_1(e_2) \mapsto e'_1(e_2)} \text{ D-APP1} \quad \frac{e_2 \mapsto e'_2}{e_1(e_2) \mapsto e_1(e'_2)} \text{ D-APP2}$$

$$\begin{array}{c}
\frac{s \cdot s', e_1 \xrightarrow{P} v_0, s' \quad s \cdot \varepsilon, e_2 \xrightarrow{P} v, \varepsilon}{s \cdot s', e_1 \text{ asa } e_2 \xrightarrow{P} v, s'} \text{ DP-ASA} \\
\frac{s \cdot s', e_1 \xrightarrow{P} v, s' \quad [v/\$!]e_2 \rightarrow \text{true}}{s \cdot s', e_1 \text{ where } e_2 \xrightarrow{P} v, s'} \text{ DP-WHERETRUE} \quad \frac{s \cdot s', e_1 \xrightarrow{P} v, s' \quad [v/\$!]e_2 \rightarrow \text{false}}{s \cdot s', e_1 \text{ where } e_2 \xrightarrow{P} \text{error}, s \cdot s'} \text{ DP-WHEREFALSE} \\
\frac{}{\varepsilon \cdot s', \text{structtcurts} \xrightarrow{P} \langle \rangle, s'} \text{ DP-STRUCTBASE} \quad \frac{}{s \cdot s', \text{structts}, \text{field}^0, \text{field}^i : e} \text{ DP-STRUCTSTEP}
\end{array}$$

8 Choosing An Appropriate Representation

We need to account for a fundamental optimization—that which allows us to choose an appropriate interpretation in the target language for a specific layout. For instance, we want to interpret a `uint(8) list` in C as a `struct {int len; char *arr};` instead of as a `struct {int len; BIGNUM *arr};`.

Let us consider a term $t : \tau$ layout. The parser which corresponds to t will parse only a subset of values of type τ . For instance, a value parsed by a parser for the term `uint(8)` must be an integer between 0 and 255 inclusively. If $t : \tau$ layout, let τ_t denote the subset of the values of type τ that t might parse.

Our requirement for the embedding type e in the target language for the layout term t is that $\tau_t \subseteq e \subseteq \tau$ or equivalently $\tau_t \leq e \leq \tau$. The relation between e and τ is witnessed by coercions: $c_{e \rightarrow \tau}$ and $c_{\tau \rightarrow e}$. We do not require that $e = \tau_t$ since τ_t may not be representable in the target language, or even if encodable, τ_t may not be an efficient representation. Furthermore, calculating a “smallest” e in a particular language’s type system is often NP-Hard (see Appendix I).

A first pass approximation for finding reasonably appropriate representations is to use the range information that is easily provided by the type constructors themselves. The use is straight forward when our integer type constructors are applied to constants as in the case of `uint(8)`. We know immediately that any integral type that can represent the range $[0 \dots 255]$ will do. A slightly more difficult example is that of the SWF RECT, but this still seems within reason: we calculate the required range of the field $xmin$ by $\bigvee_{i=1}^{31} \text{uint}(i)$ to find that any integral type that can represent $[0 \dots 2, 147, 483, 647]$ will do.

Unfortunately, we might not always want to perform that calculation. For instance, we would not want to try the same trick with the field x in the specification:

```

WOAH =
struct
    n : uint(16),
    m : uint(16),
    p : uint(16),
    x : uint(5*n*n*m - p*n + p*p - n*n*p),
tcurts

```

XXX Insert intuitions here that we will do simple range analysis on BLT Specifications to choose appropriate types. Abstract Interpretation á la Cousot and Cousot Anyone?

9 Distinguishability

An important goal not covered heretofore is that of warning the specification writer that he may have written an ambiguous grammar. We prove that determining whether a BLT specification is ambiguous is NP Hard in Appendix II.

Despite the fact that this problem is NP Hard, there are good strategies for making unambiguous grammars. The most straight-forward approach is *tagging*: guaranteeing that a bit pattern specified by one branch of a union will not appear in another branch.

Deriving this information requires the interpretation of where-clauses by our primitive layout constructors. With the base constructors we have given, we are guaranteed that singleton values will correspond to unique bit-patterns. This need not be the case in general, though, as the following example illustrates:

(if $x > 0$ then uint(8) else uint(32)) where $! = 5$

XXX Display an algorithm we will use to attack this problem, explain the subtleties of lists and follow sets, as well as unions.

Appendix I : NP Hardness of BLT Embedding Optimization Problem

[BLT-OPT] BLT Optimal Embedding Problem

INSTANCE: BLT Specification S , Lattice of Integer Types (i.e. ranges) in a target language such as C

QUESTION: What is the optimal embedding for each named layout term $t \in S$ as a type constructed using the given integral types.

We will show that MAX2SAT is reducible to BLT-OPT.

Consider the lattice of integral types one might consider for an embedding in C: unsigned char (8 bit), unsigned short (16 bit), unsigned int (32 bit), unsigned long (64 bit), and BIGNUM. Suppose Opt is a function which calculates the optimal embedding $e \geq t_\tau$ given this lattice for a layout term t in a BLT Specification.

Let $v_1, v_2 \dots v_n, \{x_i, y_i\}_{i=1}^m, k$ be an instance of the MAX2SAT over n variables and m clauses. Here each x and y is either v_j or $\neg v_j$. Let $\hat{x} \equiv 1 - v_j$ if $x = \neg v_j$ and $\hat{x} = v_j$ if $x = v_j$.

Examine the following BLT Specification:

```
ZERO_OR_ONE = uint(8) where  $! = 0$  OR  $! = 1$ 
GT_K_SAT =
struct
  v1 : ZERO_OR_ONE,
  v2 : ZERO_OR_ONE,
  ⋮
  vn : ZERO_OR_ONE,
  z : uint(16) where if  $\sum_{i=1}^m \hat{x}_i * \hat{y}_i > k$  then  $! = 255$  else  $! = 65535$ 
tcurts
```

Then $Opt(GT_K_SAT) = \text{struct}\{\text{char } x_1; \text{char } x_2; \dots \text{char } x_n; \text{unsigned short } z;\}$; if and only if the MAX2SAT instance is satisfiable.

Appendix II : NP Hardness of Ambiguity Detection

[BLT-AMBIGUITY] BLT Ambiguity Problem

INSTANCE: BLT Specification S

QUESTION: Can a sequence of bits be interpreted differently if we were to construct a non-predictive parser. (We could create a non-predictive parser by choosing non-deterministically the order of branches tried in unions and the stopping points for lists of unknown length.)

We will show that MAX2SAT is reducible to BLT-AMBIGUITY. Suppose Amb can decide whether a BLT specification is ambiguous. Examine the following extension of the BLT Specification in Appendix I:

```
ALWAYS_YES =
struct
  v1 : ZERO_OR_ONE,
  v2 : ZERO_OR_ONE,
  ⋮
  vn : ZERO_OR_ONE,
  z : uint(16) where $! = 65535
tcurts
-----
AMBIGUOUS =
union
  K_CLAUSES_SATISFIABLE : GT_K_SAT
  K_CLAUSES_UNSATISFIABLE : ALWAYS_YES
noinu
```

Then *Amb*(*AMBIGUOUS*) is true if and only if the MAX2SAT instance is satisfiable.