

# Symbolic Pointer Analysis

Jianwen Zhu

Department of Electrical and Computer Engineering  
University of Toronto, Ontario M5S 3G4, Canada  
jzhu@eecg.toronto.edu

## Abstract

One of the bottlenecks in the recent movement of hardware synthesis from behavioral C programs is the difficulty in reasoning about runtime pointer values at compile time. The pointer analysis problem has been investigated in the compiler community for two decades, which has yielded efficient, polynomial time algorithms for context-insensitive analysis. However, at the accuracy level for which hardware synthesis is desired, namely context and flow sensitive analysis, the time and space complexity of the best algorithms reported grow exponentially with program size. In this paper, we present our first step towards a new analysis technology which potentially leads to almost-linear time complexity and sub-linear space complexity algorithm even for the most accurate analysis. The key idea that contributes to this efficiency is to implicitly encode the pointer-to relation in the Boolean domain, thereby capturing the procedure transfer function completely, compactly and canonically. This represents a wide departure from the traditional techniques, all of which explicitly capture pointer-to relation using variations of point-to graph, which have to be re-evaluated for different calling contexts. Experiments for our first flow-insensitive algorithm on common benchmarks show promising result.

## 1 Introduction

Today's system-on-chips are heterogeneous: while the majority of the functionality can be handled by software running on on-chip microprocessors, often times it is the set of on-chip dedicated cores that defines the performance, and therefore the competitiveness of the chip. The design of such cores, however, is known to be time consuming under the current commercial design methodology. The effort to automatically synthesize hardware from algorithm-

level specification, known as behavioral (high-level) synthesis, has rarely gone beyond programs with simplistic data structures, such as scalars and arrays. This constraint seriously limits the applicability of behavioral synthesis technology, since any interesting applications will involve the use of complex data structures. This problem has become more apparent in the recent movement of behavioral synthesis from C/C++ programs, which encourage the use of pointers.

The presence of data pointers imposes challenges for temporal optimization, since flow-dependency between operations have to be conservatively approximated whenever pointer dereferences are involved, as well as spatial optimization, since the life times of memory blocks can not be accurately determined, and therefore their space cannot be aggressively shared. The presence of function pointer makes it difficult for whole-program analysis and therefore inter-procedural optimization, since even the call graph cannot be readily constructed. The degree at which these aggressive optimization can be carried out for pointer-intensive programs largely depends on the how well one can approximate the runtime pointer values at compile-time, which unfortunately, is an undecidable problem.

Luckily, the synthesis research community is not alone for the so-called *pointer analysis* problem, or in a slightly different formulation, the *alias analysis* problem. In fact, the optimizing compiler and programming language community has attacked this problem for two decades and a rich set of results have been yielded in the literature. The reported analysis algorithms vary with different accuracy speed tradeoff and can be categorized according to two criteria: *flow sensitivity* and *context sensitivity*. A flow-insensitive (FI) algorithm ignores the order of statements when it calculates pointer information, whereas a flow-sensitive (FS) algorithm takes control flow within a procedure into account. A context-insensitive (CI) algorithm does not distinguish the different calling contexts of a procedure, whereas a context-sensitive (CS) does.

Fast polynomial algorithms, such as derivatives of Steensgaard's [1] and Andersen's [2], have been developed for context-insensitive analysis. However, the accuracy of such analysis, especially Steensgaard's, is not acceptable for hard-

ware synthesis aiming for aggressive optimization. On the other hand, the context-sensitive analysis algorithms developed in the literature suffer from a worst-case exponential time complexity. The fundamental reasons for such complexity is three fold: First, these algorithms cannot obtain the *complete transfer function* (CTF) of a procedure, in other words, summarize the effect of procedure on program state with respect to arbitrary initial state. Therefore, each procedure has to be re-analyzed for each calling context. This problem has been alleviated by the landmark work of Wilson and Lam [3], who developed the concept of *partial transfer function*. Nevertheless, this in-efficiency is not completely eliminated. Second, the cost of propagating the effect of callee to the caller is high, since the point-to relation of the caller has to be duplicated at each call site. This cost remains to be high even with the sparse representation [4]. Third, the data structure that these algorithms use to capture program state are inherently redundant, thereby causing in-efficiency in space. For flow-sensitive analysis, where program state for each program point has to be kept, the space in-efficiency will eventually limit the size of the program one can analyze.

In this paper, we describe an analysis technology, called symbolic pointer analysis, which starts a new direction in combating the three sources of in-efficiencies. The contributions of this paper is as follows: First, we map each program memory location of interest into a cube in a Boolean space, called the *domain space*, and a cube in another Boolean space, called the *range space*. This allows us to *implicitly* represent sets of memory locations using Boolean functions, thereby providing space efficiency and potentially speed efficiency. Second, instead of using set of alias pairs, or point-to graph, we propose a new approach to represent the program state using a Boolean function in the combined domain and range space. This allows us to perform state query using implicit image computation originally developed for reachability analysis in formal verification [5, 6], rather than the explicit graph traversal. Third, we introduce the *initial state blocks*, a concept similar but different from Wilson’s extended parameters, as well as the concept of *scalar predicate*, which is necessary if accuracy due to “strong update” is not to be compromised. These two concepts allow us to capture the transfer function of a procedure *completely* as a Boolean function, which was otherwise impossible. Forth, we engineer a CSFI analysis algorithm based on the above concepts using Bryant’s Binary Decision Diagram (BDD) [7] and illustrate why the compactness and canonicity of BDD can be exploited for efficiency.

The rest of the paper includes a detailed description of the proposed representation (Section 3) and the fundamental algorithms (Section 4) under certain simplifying assumptions, which will later be relaxed. The engineering side of the algorithm is discussed in Section 5. Lastly, while we are not able to perform a comprehensive and thorough empirical comparison with other works at this point due to the time and resource constraint, we present our preliminary result

on standard benchmarks, which promisingly confirmed the expected efficiency.

In the text that follows, we use the *formal algorithm notation* (FAN) to state definitions and describe algorithms. Unlike pseudo-code based algorithm description, FAN relies on a type system, where each type is represented by a set, to present the algorithm in a formal, precise manner. For example, we use the notation  $\langle \rangle^A$  to represent the power set of  $A$ , therefore any value of type  $\langle \rangle^A$  must be a set of values of type  $A$ . Similarly, we use the notation  $[ ]^A$  to represent the set of all sequences over elements of  $A$ , therefore any value of type  $[ ]^A$  represents a sequence of values of type  $A$ . Readers are expected to find this notation very similar to any strongly-typed programming languages and hence straightforward to be translated into implementation, yet abstract enough to allow concise presentation.

## 2 Related Work

To quote a recent excellent survey by Hind [8], over seventy-five papers and nine PhD thesis haven be published on pointer analysis. So here we only review important milestones and recent works that are relevant to ours.

In the category of FICI pointer analysis, Steensgaard’s work [1] stands out as the first equality-based method, which treats assignment as bidirectional and use a union-find data structure. His analyzer is extremely fast and has analyzed million lines of industrial code, however, the precision of equality-based approach degrades very fast in general, even with later improvement [9]. Andersen’s [2] popular subset-based improves precision by treating assignment as a unidirectional flow of values, yet still run’s in polynomial time.

Our algorithm primarily works in the research space of accurate FS/CS analysis. Wilson and Lam [4] proposed the concept of partial transfer function in an effort to reduce the number of times a procedure has to be analyzed. Along the same line, recently Liang and Harrold [10] has proposed to use the local and global auxiliary blocks, similar to our initial state block. However, their model cannot extend beyond FS analysis since it is impossible to detect strong updates.

The application of pointer analysis algorithm for behavioral synthesis is emerging. Semeria and De Micheli have researched on how to encode C pointers in hardware [11]. By a combination of runtime test and compile time analysis, they can produce efficient, customized pointer implementation. Panda et. al. [12] applied pointer analysis to optimally determine array memory layout for cache performance. Zhu [13] applied pointer analysis to the problem of memory block minimization.

## 3 Abstracting Program State

The goal of pointer analysis is to statically estimate the runtime program state, or the set of values each program location can hold. To trade accuracy for analysis speed, we often

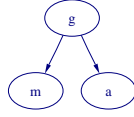


Figure 1: Program state on completion of program in Example 1.

collapse related program locations together, thereby forming a **block**. Locations within a block are not distinguished. The blocks can be global variables, or local variables, or procedure parameters, or dynamically allocated memory blocks. The values of interest are only the addresses of the blocks.

**Example 1** A C program modified from [10]. Note that the program contain global blocks *g*, *a*, local blocks *p*, *q*, *r*, *t*, *f*, and dynamic block allocated at *S3*.

<i>char *g, a;</i>	1	<i>char *g, a;</i>	1
<i>void main() {</i>	2	<i>void main() {</i>	2
<i>char *p, *q;</i>	3		3
	4		4
<i>S1: alloc( &amp;p );</i>	5	<i>call alloc, [&amp;p];</i>	5
<i>getg( &amp;q );</i>	6	<i>call getg, [&amp;q, 0];</i>	6
<i>g = &amp;a;</i>	7	<i>store g, &amp;a;</i>	7
<i>}</i>	8	<i>}</i>	8
	9		9
<i>void getg( char** r ) {</i>	10	<i>getg( [r1, g1] ) {</i>	10
<i>char **t = &amp;g;</i>	11	<i>store t, &amp;g;</i>	11
<i>if( g == NULL )</i>	12	<i>others</i>	12
<i>S2: alloc( t );</i>	13	<i>call alloc [t];</i>	13
<i>*r = *t;</i>	14	<i>store *r, *t;</i>	14
<i>}</i>	15	<i>}</i>	15
	16		16
<i>void alloc( char** f ) {</i>	17	<i>alloc( [f1] ) {</i>	17
<i>S3: *f = malloc(1);</i>	18	<i>store *f, &amp;m;</i>	18
<i>}</i>	19	<i>}</i>	19

(a) C source code

(b) preprocessed

The program state is often abstracted as a *point-to graph*, whose vertices represent the set of blocks, and an edge  $\langle u, v \rangle$  from block *u* to block *v* indicates that it is possible that the content of block *u* is the address of block *v*. The set of all edges defines the point-to relation.

**Example 2** Figure 1 shows a point-to graph capturing the program state after the completion of the main procedure in Example 1. Note here the block *m* represents the dynamically allocated block in *alloc*.

The out degree of a vertex is potentially unbounded, as in Andersen's algorithm [2]. The size of the point-to graph poses a serious problem, especially for FI analysis where a program state has to be maintained for each program point. Compromise has been made so that the out degree is K-limited, as in Horowitz's algorithm [14], or 1-limited, as in Steensgaard's algorithm [1]. For the latter two, further collapsing of blocks into larger blocks is necessary and precision degrades quite rapidly during the process.

We now propose an alternative way to capture the point-to relation. As shown in Definition 1, we associate each block with two Boolean functions, called its domain and range. The set of domain functions of all blocks form a orthogonal function set in the Boolean space, called the domain space, spanned by the set of Boolean variables  $X^* = \{x_1^*, x_2^*, \dots, x_n^*\}$ , such that  $\forall u \neq v, u.domain \times v.domain = 0$ . Similarly, the set of range functions of all blocks form a orthogonal function set in the companion Boolean space, called the range space, spanned by the set of Boolean variables  $X = \{x_1, x_2, \dots, x_n\}$ .

**Definition 1** A block is a member of

<i>SpaBlock = tuple {</i>	1
<i>kind</i>	2
<i>domain</i>	3
<i>range</i>	4
<i>}</i>	5

where *SpaDD* is the set of Boolean functions.

The domain and range functions of blocks are most conveniently selected as disjoint cubes in the Boolean spaces. Example 3 shows the assignment of cubes to the blocks.

**Example 3** Cube assignment for blocks in Example 1.

Block	domain	range
a	$\bar{x}_1 \bar{x}_2 \bar{x}_3$	$\bar{x}_1 \bar{x}_2 \bar{x}_3$
g	$\bar{x}_1 \bar{x}_2 x_3$	$\bar{x}_1 \bar{x}_2 x_3$
p	$\bar{x}_1 x_2 \bar{x}_3$	$\bar{x}_1 x_2 \bar{x}_3$
q	$\bar{x}_1 x_2 x_3$	$\bar{x}_1 x_2 x_3$
t	$x_1 \bar{x}_2 \bar{x}_3$	$x_1 \bar{x}_2 \bar{x}_3$
r	$x_1 \bar{x}_2 x_3$	$x_1 \bar{x}_2 x_3$
f	$x_1 x_2 \bar{x}_3$	$x_1 x_2 \bar{x}_3$
m	$x_1 x_2 x_3$	$x_1 x_2 x_3$

We can now capture the point-to relation by mapping each edge  $\langle u, v \rangle$  in the point-to graph by a Boolean product  $u.domain \times v.range$ . In other words, given a program state represented by *E*, it can be represented by a Boolean function  $\sum_{\langle u, v \rangle \in E} u.domain \times v.range$ .

**Example 4** The program state in Example 2 can be represented by a Boolean function:  $\bar{x}_1 \bar{x}_2 x_3 x_1 x_2 x_3 + \bar{x}_1 \bar{x}_2 x_3 \bar{x}_1 \bar{x}_2 \bar{x}_3$ .

The point-to graph can be used to represent the program state in a procedure only when the program state before the procedure is called is known. To safely capture the point-to relation under all circumstances, one either resolve to the CI approach where program states at all call sites are merged, thereby sacrificing precision, or try to drive the so-called transfer function, which can be intuitively considered as point-to relation parameterized over different calling context. The parameters of the transfer function are essentially the *content* of all memory blocks that are accessed within the procedure,

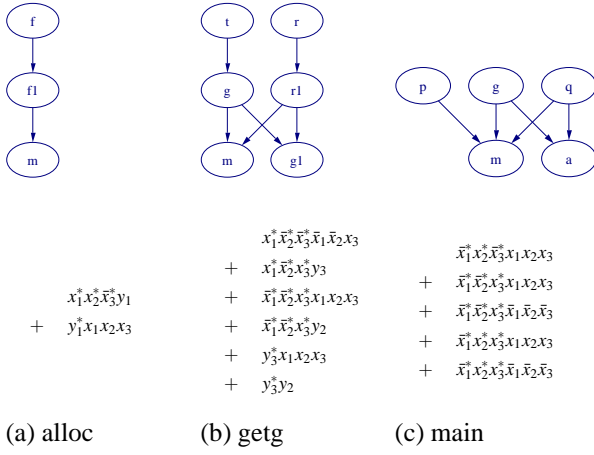


Figure 2: Transfer functions.

or the procedures directly or indirectly called by the procedure. One can treat a parameter as if it is a separate block, called the **initial state block**.

**Example 5** Consider the procedure *alloc* in Example 1, where the parameter *f* is dereferenced. Since the value of *f* is not known, what memory block will be updated after statement *S3* cannot be determined. With the introduction of the initial state block *fl*, the procedure can be summarized with a transfer function as shown in the point-to graph of Figure 2 (a).

Again, we use disjoint cubes to encode initial state blocks in the initial state domain space spanned by  $Y^* = \{y_1^*, \dots, y_m^*\}$  and range space spanned by  $Y = \{y_1, \dots, y_m\}$ . Example 6 shows one possible assignment.

**Example 6** cube assignment for initial state blocks in Example 1.

Block	domain cube	range cube
fl	$y_1^*$	$y_1$
gl	$y_2^*$	$y_2$
rl	$y_3^*$	$y_3$

Given this assignment, the transfer functions of the procedures are shown in Figure 2. Since local variables and parameters of a procedure are not visible outside their scope, the related point-to relation shall not be propagated to its caller. So often, the transfer function is pruned, as shown in Figure 3.

## 4 Symbolic Pointer Analysis Algorithm

### 4.1 Program Modeling

We now describe our pointer analysis algorithm. In order to focus on the fundamentals, rather than the implementation

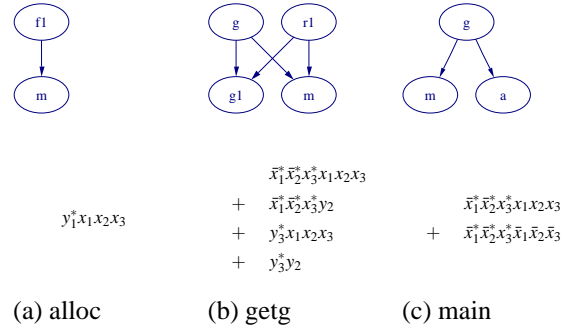


Figure 3: Transfer functions after pruning.

details, we made a number of simplifying assumptions for the ease of presentation.

We assume a program is preprocessed so that it is available in the form defined in Definition 2, which contains a set of blocks. The blocks can be globals, locals, parameters, initial blocks and procedures, as discussed in Section 3. All blocks are treated as scalars. The additional information for procedure blocks can be obtained from *procs*. The program also contains the Boolean environment necessary to define symbolic point-to relation, as discussed in Section 3. Since a block correspond to a cube in the Boolean space, a Boolean function can be used to represent a set of blocks. The utility function *spaEnumerate* can be used to enumerate all blocks corresponding to a Boolean function.

**Definition 2** A program is a member of

<i>Spa</i> = tuple {	6
<i>blks</i> : $\langle \rangle \text{SpaBlock};$	7
<i>procs</i> : $\text{SpaBlock} \mapsto \text{SpaProc};$	8
<i>domainVars</i> : $\langle \rangle \text{SpaDD};$	9
<i>rangeVars</i> : $\langle \rangle \text{SpaDD};$	10
<i>cache</i> : $\text{SpaDD} \times \text{SpaDD} \times \{\text{bddAnd}, \text{bddOr}, \dots\} \mapsto \text{SpaDD};$	11
}	12
<i>spaEnumerate</i> = func( <i>spa</i> : <i>Spa</i> , <i>l</i> : <i>SpaDD</i> ) : $\langle \rangle \text{SpaBlock};$	13
	14

A procedure contains a set of parameter, local, heap allocated blocks, and a series of instructions. Note that the formals contains a sequence of initial state blocks which are originally implicit in the source program.

**Definition 3** A procedure is a member of

<i>SpaProc</i> = tuple {	15
<i>formals</i> : $\langle \rangle \text{SpaBlock};$	16
<i>params</i> : $\langle \rangle \text{SpaBlock};$	17
<i>locals</i> : $\langle \rangle \text{SpaBlock};$	18
<i>heaps</i> : $\langle \rangle \text{SpaBlock};$	19
<i>instrns</i> : $\langle \rangle \text{SpaInstrn};$	20
}	21

As shown in Figure 4, we consider only two types of instructions, assignment instruction (store), call instruction (call), and ignore others. The destination and sources of the instruction are in the form of access path.

**Definition 4** An instruction is a member of

<i>SpaInstrn</i> = tuple {	22
<i>opcode</i> : {store, call, others};	23
<i>dst</i> : <i>SpaAccessPath</i> ;	24
<i>srcs</i> : [] <i>SpaAccessPath</i> ;	25
}	26

We assume an access path can be characterized by a block, which indicates the root of access, and level, which indicates the number of dereferences performed on the block. The C expression  $\&b$  corresponds to level -1,  $b$  corresponds to level 0, and  $***b$  corresponds to level 3.

**Definition 5** An access path is a member of

$SpaAccessPath = \text{tuple} \{$	27	
$blk$	$: SpaBlock;$	28
$level$	$: \mathbb{Z};$	29
$\}$		30

We assume that the procedure does not contain recursive calls. Also we do not distinguish memory allocated at the same site under different calling contexts.

The program in Example 1 (a) is preprocessed into Example 1 (b). Note that an important difference from its source is that the formals refers to the initial state blocks, which are not explicit in the source code before. The corresponding actuals at the callsite, are also made explicit.

## 4.2 Pointer Analysis Algorithm

Algorithm 1 computes the Boolean CTF for a given procedure  $root$ . The parameter  $tfs$  keeps track of the CTFs of the procedures that have been analyzed. The algorithm starts by computing the initial points-to relation, which involves linking the blocks referenced in the procedures with their initial state blocks present in the formals. It then evaluates each instruction, after which it updates the current state  $state$ . Since this is a FI algorithm, the order of the instruction does not matter. When it encounters a store instruction, it calls  $spaUpdateState$  to compute the new state  $nstate$ . When a call instruction is encountered, the algorithm first calls  $spaQueryState$  to find out the set of procedures that the call instruction may call. This process is necessary since we may be accessing function pointers. The algorithm then iterate through all the potential callees, and call  $spaApply$  to apply the callsite point-to information to the CTF of the callee, thereby propagating the point-to information of the callee

back to the call site. If the CTF is not available at  $tfs$ , it recurs and obtain the CTF of the callee before proceed. This process repeats until a fixed-point is reached, when there is no further update of the state function. Note that while this is an iterative process, each procedure needs to be analyzed only once. Before it returns the result, a call to  $spaPrune$  will be performed in order to remove from CTF the set of point-to edges emitting from local and parameter blocks, which, by the language scoping rule, are not accessible outside.

**Algorithm 1** Pointer analysis.

```

spaAnalysisFICS = func(                                     31
  spa : Spa, root : SpaBlock, tfs : SpaBlock  $\mapsto$  SpaDD    32
) : SpaDD {                                               33
  var state, nstate : SpaDD;                               34
  var changed : {true, false};                             35
  changed = true, state = spaInitState( spa, root );       36
  while( changed == true ) {                               37
    changed = false;                                       38
    forall( instrn  $\in$  spa.procs(root).instrns ) {         39
      if( instrn.opcode == store ) {                       40
        nstate = spaUpdateState(spa, state, instrn.dst, instrn.srcs[0]); 41
        if( nstate  $\neq$  state )                             42
          changed = true, state = nstate;                 43
      }                                                     44
    } else if( instrn.opcode == call ) {                   45
      forall( callee  $\in$  spaEnumerate(                       46
        spaQueryState( spa, state, instrn.dst.blk.range, instrn.dst.level ) 47
      ) ) {                                                 48
        if( tfs(callee) ==  $\emptyset$  )                       49
          spaAnalysisFICS(spa, callee, tfs);               50
        nstate = spaApply(spa, state, instrn.srcs, callee, tfs(callee)); 51
        if( nstate  $\neq$  state )                             52
          changed = true, state = nstate;                 53
      }                                                     54
    }                                                       55
  }                                                         56
  state = spaPrune( state, root.params  $\cup$  root.locals ); 57
  tfs(root) = state;                                       58
  return state ;                                          59
}                                                         60

```

The following results ensures that the algorithm terminates.

**Lemma 1** A set of Boolean functions defined on the Boolean space  $B^n$  forms a lattice, where the bottom element is constant 0 function, the top element is constant 1 function, and the meet operator is the Boolean or operator.

**Lemma 2** The program state in the form of a Boolean function computed by Algorithm 1 at each iteration is monotonically increasing with respect to the Boolean or operator.

**Theorem 1** Algorithm 1 terminates.

## 4.3 State Query

Algorithm 2 is used to find out the set of all possible values an access path may have under a program state. It is needed to process every access path, be it a destination or source, of an instruction. On the point-to graph, it is equivalent to

determining the envelop of the  $n$ th level breath-first search starting from a block, when  $n$  is the level of the access path. In our Boolean framework, the algorithm is similar to the image computation algorithm in reachability analysis. For level 1 access  $*b$ , where  $b$  is a block, one only needs to multiply the state function with the domain cube of the block. The resultant Boolean function captures the set of blocks in the range space. To access one more level, one needs to call *bddMirror* to isomorphically map the computed range back to the domain space, which can again be used to multiply with the state function. This process continues until we reach the desired level. Note that unlike algorithms based on point-to graph, the image computation can be performed implicitly, in other words, it can be performed purely logic operation without resolving to explicit graph traversal.

#### Algorithm 2 State query.

```

spaQueryState = func(                                     63
  spa : Spa, state : SpaDD, from : SpaDD, level : Z      64
) : SpaDD {                                              65
  if( level == 0 ) return from ;                          66
  return bddAnd(                                          67
    spa, state, bddMirror(                                68
      spa, spaQueryState(spa, state, from, level-1)      69
    )                                                       70
  );                                                       71
}                                                         72

```

### 4.4 Evaluating Stores

We use Algorithm 3 to compute the new state whenever a store instruction is encountered. Consider a simple case where the access paths for both the destination and the source of the instruction are simple blocks, which correspond to adding an edge in point-to graph based algorithms. Here we simply need to add a product term to the current state, which is the product of the domain cube of the destination and the range cube of the source. The operation holds equally well with complex access paths, whose net effect will still amount to adding a product of two Boolean functions. Note that this corresponds to collectively adding many point-to edges at the same time!

#### Algorithm 3 State update.

```

spaUpdateState = func(                                     73
  spa : Spa, state : SpaDD, dst : SpaAccessPath, src : SpaAccessPath 74
) : SpaDD {                                              75
  return bddOr(                                          76
    state, bddAnd(                                       77
      bddMirror( spa,                                     78
        spaQueryState( spa, state, dst.blk.range, dst.level ) 79
      ),                                                  80
      spaQueryState( spa, state, src.blk.range, dst.level+1 ) 81
    )                                                       82
  );                                                       83
}                                                         84

```

### 4.5 Evaluating Calls

The pointer information of the callee has to be propagated back to the caller by applying the program state at the call site to the transfer function of the callee. This is achieved by Boolean function composition, in other words, substituting the Boolean variables that construct the initial state blocks by a set of Boolean functions. Such projection map can be obtained from the binding of actuals to the formal. If we use a single variable for each initial state block, it can be computed trivially.

#### Algorithm 4 Evaluate call.

```

spaApply = func(                                           85
  spa : Spa, state : SpaDD, srcs : [ ]SpaAccessPath,      86
  proc : SpaBlock, tf : SpaDD                             87
) : SpaDD {                                              88
  var bind, proj : SpaDD  $\mapsto$  SpaDD;                  89
  var actual, formal : SpaDD;                             90
  forall( i  $\in$  {0..|srcs|-1} ) {                          91
    actual = spaQueryState(spa, state, srcs[i].blk.range, srcs[i].level+1 ); 92
    formal = proc.formals[i].range;                       93
    bind(formal) = actual;                                 94
    bind(bddMirror(spa, formal)) = bddMirror(spa, actual); 95
    proj = bddFindProjection(spa, bind);                  96
  }                                                       97
  return bddOr(                                           98
    state, bddCompose( spa, tf, proj )                   99
  );                                                       100
}                                                         101

```

### 4.6 State Pruning

Algorithm 5 first computes the Boolean function that corresponds to the set of all blocks need to be pruned. It then predicates the state with the inversion of the computed function.

#### Algorithm 5 State pruning.

```

spaPrune = func(                                           103
  spa : Spa, state : SpaDD, vars : [ ]SpaBlock           104
) : SpaDD {                                              105
  var sum : SpaDD;                                       106
  sum = bddFalse(spa);                                   107
  forall( v  $\in$  vars )                                     108
    sum = bddOr( sum, v.domain );                       109
  return bddAnd( spa, state, bddNot( spa, sum ) );       110
}                                                         111

```

### 4.7 Relaxing Assumptions

We now discuss how the basic algorithms described can be modified if the simplifying assumptions are relaxed.

C extensively uses aggregates such as arrays and structs. The notions of block as well as access path need to be refined so that different fields of a struct can be distinguished. Wilson [4] developed the concept of a *location set*, which

characterizes a set of grouped location using a pair of normalized integers, called offset and stride. The use of location set won't change our algorithm, however, in order to maintain the invariant that blocks are disjoint, some block collapsing needs to be performed as the algorithm proceeds. Note that this tradeoff is also made in [4].

The basic algorithm does not handle *strong update*, which concerns the assignment of a unique location. Global and local scalar blocks are always unique. Dynamic blocks are never unique. When a strong update happens, instead of merging the new value with the original value as the new point-to set, as in the case of *weak update*, the original value stored in the location under current state needs to be removed. Under our symbolic framework, this removal can be easily carried out the same way as Algorithm 5. The real challenge, however, is that for initial state blocks, the knowledge of whether they are unique scalars is not available. This factor alone leads to the failure of true CTF in previous efforts [4, 15]. On the other hand, under our symbolic framework, it is just as simple as introducing an extra Boolean variable, called the *scalar predicate*, for each initial block.

Some extra care needs to be exercised to handle programs with recursive functions. We follow the usual path: whenever a calling cycle is hit, we use the transfer function computed so far as an approximation, and will iterate until a fixed-point is reached.

## 5 Engineering a Fast Algorithm

In addition to many implementation details that we choose to omit in the algorithm presentation, one extremely important issue left undiscussed is how we actually represent and manipulate the Boolean functions which forms the foundation of our algorithm. A poor choice of the Boolean function representation may invalidate the entire methodology. Our use of Bryant's Binary Decision Diagram (BDD) [7] is rather an obvious decision. It is a canonical, compact, and high-quality manipulation package is freely available [16].

Some considerations during the construction of BDD prove to be extremely helpful for the algorithm efficiency. One of the most commonly used BDD operation is *bddMirror*, which maps a Boolean function in range space to an equivalent one in domain space. We keep this operation cheap by aligning the variable ordering between the two variable sets: while the variable order within a set can be dynamically changed to reduce overall BDD size, they keep the same order across two different set.

An interesting argument for discussion is that since BDD is a graph-based representation of a Boolean function, why don't we use the point-to graph in the first place, which seems to be much more intuitive? Our observations are as follows: First of all, point-to graphs need to be maintained for every procedure, each of which may share many common edges. In other words, there are a large amount of

redundancy. In contrast, BDD enables the maximum sharing among graph nodes, and point-to information in different procedures, at different program points can be reused. This space efficiency becomes crucial for the accurate analysis of large programs.

Second, many otherwise expensive operations can be performed implicitly using BDD. The complexity of an operation now depends on the size of the BDD, rather than the size of the point-to graph. Just like image computation, operations in BDD can be batch processed on a large collection of blocks.

Third, an often overlooked factor is the canonicity of BDD representation. The ease of defining a signature on an operation applied on two BDDs allows us to cache and reuse the result easily and on a large scale. This efficiency is in essence the same as the principle of dynamic programming: if a subproblem can be uniquely identified, it should be solved only once and its result should be shared by other upper-level problems.

## 6 Experimental Result

We implemented an FICS pointer analysis tool based on the ideas outlined in this paper using a synthesis/compiler infrastructure we developed. Our infrastructure includes several frontends such as C, Java, Verilog as well as a C-based, objected oriented, polymorphic language developed specifically for System-on-chip. We also leverage Somenzi's CUDD package [16] for BDD manipulation.

Our first implementation wasn't intended for comprehensive comparison with what has been published in the literature, but rather a validation of a new idea. We have important omissions: e.g., we do not yet distinguish field elements, and have not fully implement Wilson's location set concept. Our transfer functions for libraries are hardwired, and we still do not have a strategy to summarize library function without analyzing the library source code. We have not yet handle non-local control transfers such as *setjmp/longjmp*. All these complications have prevented us from attempting large industrial benchmarks. Also we are not even close to reaching the last miles of performance improvement using symbolic technology. Therefore, at this stage we are not ready for a fair, direct performance comparison with traditional techniques.

Nevertheless, we were able to apply our implementation on the benchmark set commonly used in the pointer analysis research community: one set from McGill university [17], and the other from Dr. Landi [18]. The statistics of the benchmarks are shown in Table 1, where the third column indicates the lines of code, the fourth column indicates the number of procedures (not including libraries), and the fifth column indicates the percentage of control flow nodes that have pointer assignment. The last column shows the runtime of our algorithm obtained on a Ultra-SPARC 5. The result is very encouraging: all experiments finish in seconds.

Name	Source	LOC	#procs	density	Run Time (s)
01.qbsort	McGill	325	8	24.1%	0.10
06.matx	McGill	350	7	13.5%	0.13
15.trie	McGill	358	13	23.4%	0.21
04.bisect	McGill	463	9	9.7%	0.10
17.bintr	McGill	496	17	8.8%	0.13
05.eks	McGill	1202	30	4.0%	0.20
08.main	McGill	1206	41	20.9%	1.33
09.vor	McGill	1406	52	28.6%	5.54
allroots	Landi	227	7	1.3%	3.02
football	Landi	2354	58	1.8%	2.38
compiler	Landi	2360	40	5.1%	5.3
assembler	Landi	3446	52	16.6%	10.63
simulator	Landi	4639	111	6.3%	4.03

Table 1: Experimental results

## 7 Conclusion

In this paper, we described an idea for accurate pointer analysis in a way significantly different from traditional approaches. We discuss the rationale of this approach by first analyzing sources of in-efficiency in traditional methods. We demonstrated our preliminary implementation and experiment, which confirmed our expectation. In the future, we will provide a thorough implementation of a comparative framework so that solid empirical data can be obtained for different methods. We will also apply pointer analysis techniques in high-level and system-level synthesis.

## References

- [1] Bjarne Steensgaard, “Points-to analysis in almost linear time,” in *Symposium on Principles of Programming Languages*, 1996, pp. 32–41.
- [2] O. Andersen, *Program Analysis and Specialization for the C Programming Language*, Ph.D. thesis, Computer Science Department, University of Copenhagen, 1994.
- [3] R. Wilson and M. Lam, “Efficient context-sensitive pointer analysis for c programs,” in *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation*, June 1995, pp. 1–12.
- [4] Robert P. Wilson and Monica S. Lam, “Efficient context-sensitive pointer analysis for c programs,” in *SIGPLAN Conference on Programming Language Design and Implementation*, 1995, pp. 1–.
- [5] O. Coudert, C. Berthet, and J. C. Madre, “A unified framework for the formal verification of sequential circuits,” in *Proceedings of the International Conference on Computer-Aided Design*, November 1990, pp. 126–129.
- [6] O. Coudert and J. C. Madre, “The implicit set paradigm: a new approach to finite state system verification,” *Formal Methods in System Design*, vol. 2, no. 6, pp. 133–145, 2001.
- [7] Randal E. Bryant, “Graph-based algorithms for Boolean function manipulation,” *IEEE Transactions on Computers*, vol. C-35, no. 8, pp. 677–691, Aug. 1986.
- [8] Michael Hind, “Pointer analysis: Haven’t we solved this problem yet,” in *ACM SIGPLAN-SIGSOFT Workshop on Program Analysis for Software Tools and Engineering (PASTE)*, June 2001.
- [9] Sean Zhang, Barbara G. Ryder, and William Landi, “Program decomposition for pointer aliasing: A step toward practical analyses,” in *Foundations of Software Engineering*, 1996, pp. 81–92.
- [10] Donglin Liang and Mary Jean Harrold, “Efficient computation of parameterized pointer information for interprocedural analyses,” in *Static Analysis Symposium*, 2001, pp. 279–298.
- [11] L. Semeria and G. De Micheli, “Resolution, optimization, and encoding of pointer variables for the behavioral synthesis from c,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, February 2001.
- [12] P. Panda, L. Semeria, and G. De Micheli, “Cache-efficient memory layout of aggregate data structures,” in *Proceedings of the International Symposium on System Synthesis*, September 2001.
- [13] J. Zhu, “Static memory allocation by pointer analysis and coloring,” in *Design Automation and Test in Europe*, March 2001.
- [14] Susan Horwitz, Thomas W. Reps, and David Brinkley, “Interprocedural slicing using dependence graphs,” in *SIGPLAN Conference on Programming Language Design and Implementation*, 1988, pp. 35–46.
- [15] Donglin Liang, Maikel Pennings, and Mary Jean Harrold, “Extending and evaluating flow-insensitive and context-insensitive points-to analyses for java,” in *Workshop on Program Analysis For Software Tools and Engineering*, 2001, pp. 73–79.
- [16] F. Somenzi, “CUDD: Binary decision diagram package release,” <http://vlsi.Colorado.EDU/~fabio/CUDD/cuddIntro.html>, 1998.
- [17] L. Hendren et al, “Mccat compiler project,” <http://www-acaps.cs.mcgill.ca/~benadmin/benchmarks/>.
- [18] W. Landi et al, “Prolangs analysis framework,” <http://www.prolangs.rutgers.edu>.