

Extending the Theory of Arrays: memset, memcpy, and Beyond

Stephan Falke, Florian Merz, and Carsten Sinz

INSTITUTE FOR THEORETICAL COMPUTER SCIENCE (ITI)

```
struct list_node {  
    int data;  
    struct list_node *tail;  
};  
typedef struct list_node list;  
  
list *reverse(list *l) {  
    list *r = l, *p = NULL;  
    while (r != NULL) {  
        list *q = r;  
        r = r->tail;  
        q->tail = p;  
        p = q;  
    }  
    return p;  
}
```



- SMT-solvers are routinely used in program analysis:
 - Deductive program verification
 - Symbolic execution
 - Software bounded model checking
 - ...

- SMT-solvers are routinely used in program analysis:
 - Deductive program verification
 - Symbolic execution
 - Software bounded model checking
 - ...
- Prominent theory: \mathcal{T}_A (theory of arrays)
 - Model arrays/structures/objects in the program
 - Model main memory

\mathcal{T}_A : The Theory of Arrays

index terms	$t_I ::= \dots$
element terms	$t_E ::= \dots \mid \text{read}(t_A, t_I)$
array terms	$t_A ::= a \mid \text{write}(t_A, t_I, t_E)$

index terms	$t_I ::= \dots$
element terms	$t_E ::= \dots \mid \text{read}(t_A, t_I)$
array terms	$t_A ::= a \mid \text{write}(t_A, t_I, t_E)$

$$p = r \implies \text{read}(\text{write}(a, p, v), r) = v$$

$$\neg(p = r) \implies \text{read}(\text{write}(a, p, v), r) = \text{read}(a, r)$$

index terms	$t_I ::= \dots$
element terms	$t_E ::= \dots \mid \text{read}(t_A, t_I)$
array terms	$t_A ::= a \mid \text{write}(t_A, t_I, t_E)$

a write modifies the position written to ...

$$\curvearrowright \quad p = r \quad \Longrightarrow \quad \text{read}(\text{write}(a, p, v), r) = v$$

$$\neg(p = r) \quad \Longrightarrow \quad \text{read}(\text{write}(a, p, v), r) = \text{read}(a, r)$$

index terms	$t_I ::= \dots$
element terms	$t_E ::= \dots \mid \text{read}(t_A, t_I)$
array terms	$t_A ::= a \mid \text{write}(t_A, t_I, t_E)$

a write modifies the position written to ...

$\hookrightarrow p = r \implies \text{read}(\text{write}(a, p, v), r) = v$

$\hookrightarrow \neg(p = r) \implies \text{read}(\text{write}(a, p, v), r) = \text{read}(a, r)$

...and nothing else

How to model standard library functions such as `memset` and `memcpy`?

```
void *memset(void *dst, int c, size_t n);
```

```
void *memcpy(void *dst, const void *src, size_t n);
```


How to model standard library functions such as `memset` and `memcpy`?

might not be constant!

```
void *memset(void *dst, int c, size_t n);
```

might not be constant!

```
void *memcpy(void *dst, const void *src, size_t n);
```

Motivation

```
...  
memcpy(a, b, 4);  
...
```

```
...  
memcpy(a, b, 4);  
...
```

$$a_1 = \text{write}(a, 0, \text{read}(b, 0))$$

```
...  
memcpy(a, b, 4);  
...
```

$$a_1 = \text{write}(a, 0, \text{read}(b, 0))$$
$$a_2 = \text{write}(a_1, 1, \text{read}(b, 1))$$

```
...  
memcpy(a, b, 4);  
...
```

$a_1 = \text{write}(a, 0, \text{read}(b, 0))$

$a_2 = \text{write}(a_1, 1, \text{read}(b, 1))$

$a_3 = \text{write}(a_2, 2, \text{read}(b, 2))$

```
...  
memcpy(a, b, 4);  
...
```

$$a_1 = \text{write}(a, 0, \text{read}(b, 0))$$
$$a_2 = \text{write}(a_1, 1, \text{read}(b, 1))$$
$$a_3 = \text{write}(a_2, 2, \text{read}(b, 2))$$
$$a' = \text{write}(a_3, 3, \text{read}(b, 3))$$

```
...  
memcpy(a, b, 4);  
...
```

$a_1 = \text{write}(a, 0, \text{read}(b, 0))$

$a_2 = \text{write}(a_1, 1, \text{read}(b, 1))$

$a_3 = \text{write}(a_2, 2, \text{read}(b, 2))$

$a' = \text{write}(a_3, 3, \text{read}(b, 3))$

Does **not** scale well for large constants

Motivation

```
...  
memcpy(a, b, n);  
...
```


Motivation

```
...  
memcpy(a, b, n);  
...
```

???

Motivation

```
...  
memcpy(a, b, n);  
...
```

$a' = \text{copy}(a, 0, b, 0, n)$

```
...  
memcpy(a, b, n);  
...
```

$$a' = \lambda i. \text{ITE}(0 \leq i < n, \text{read}(b, i), \text{read}(a, i))$$

```
...  
memcpy(a, b, n);  
...
```

$$a' = \lambda i. \text{ITE}(0 \leq i < n, \text{read}(b, i), \text{read}(a, i))$$

\Rightarrow Extend \mathcal{T}_A by λ -terms that describe arrays

Motivation

```
...  
memset(a, v, n);  
...
```

```
...  
memset(a, v, n);  
...
```

$$a' = \lambda i. \text{ITE}(0 \leq i < n, v, \text{read}(a, i))$$

```
int i, j, n = ...;
int *a = malloc(2 * n * sizeof(int));
for (i = 0; i < n; ++i) {
    a[i] = i + 1;
}
for (j = n; j < 2 * n; ++j) {
    a[j] = 2 * j;
}
```

```
int i, j, n = ...;
int *a = malloc(2 * n * sizeof(int));
for (i = 0; i < n; ++i) {
    a[i] = i + 1;
}
for (j = n; j < 2 * n; ++j) {
    a[j] = 2 * j;
}
```

$$a' = \lambda i. \text{ITE}(0 \leq i < n, i + 1, \text{read}(a, i))$$


```
int i, j, n = ...;
int *a = malloc(2 * n * sizeof(int));
for (i = 0; i < n; ++i) {
    a[i] = i + 1;
}
for (j = n; j < 2 * n; ++j) {
    a[j] = 2 * j;
}
```

$$a' = \lambda i. \text{ITE}(0 \leq i < n, i + 1, \text{read}(a, i))$$

$$a'' = \lambda j. \text{ITE}(n \leq j < 2 * n, 2 * j, \text{read}(a', j))$$

① $\mathcal{T}_{\lambda\mathcal{A}}$: an extension of $\mathcal{T}_{\mathcal{A}}$ with λ -terms

- 1 $\mathcal{T}_{\lambda\mathcal{A}}$: an extension of $\mathcal{T}_{\mathcal{A}}$ with λ -terms
- 2 Satisfiability checking for $\mathcal{T}_{\lambda\mathcal{A}}$

$\mathcal{T}_{\lambda A}$: The Theory of Arrays with λ -Terms

index terms	$t_I ::= \dots$
element terms	$t_E ::= \dots \mid \text{read}(t_A, t_I)$
array terms	$t_A ::= a \mid \text{write}(t_A, t_I, t_E)$

$$p = r \implies \text{read}(\text{write}(a, p, v), r) = v$$

$$\neg(p = r) \implies \text{read}(\text{write}(a, p, v), r) = \text{read}(a, r)$$

$\mathcal{T}_{\lambda A}$: The Theory of Arrays with λ -Terms

index terms	$t_I ::= \dots$
element terms	$t_E ::= \dots \mid \text{read}(t_A, t_I)$
array terms	$t_A ::= a \mid \text{write}(t_A, t_I, t_E) \mid \lambda i. t_E$

$$p = r \implies \text{read}(\text{write}(a, p, v), r) = v$$

$$\neg(p = r) \implies \text{read}(\text{write}(a, p, v), r) = \text{read}(a, r)$$

$\mathcal{T}_{\lambda A}$: The Theory of Arrays with λ -Terms

index terms	$t_I ::= \dots$
element terms	$t_E ::= \dots \mid \text{read}(t_A, t_I)$
array terms	$t_A ::= a \mid \text{write}(t_A, t_I, t_E) \mid \lambda i. t_E$

$$p = r \implies \text{read}(\text{write}(a, p, v), r) = v$$

$$\neg(p = r) \implies \text{read}(\text{write}(a, p, v), r) = \text{read}(a, r)$$

$$\text{read}(\lambda i. s, r) = s[i/r]$$

index terms	$t_I ::= \dots$
element terms	$t_E ::= \dots \mid \text{read}(t_A, t_I)$
array terms	$t_A ::= a \mid \text{write}(t_A, t_I, t_E) \mid \lambda i. t_E$

$$p = r \implies \text{read}(\text{write}(a, p, v), r) = v$$

$$\neg(p = r) \implies \text{read}(\text{write}(a, p, v), r) = \text{read}(a, r)$$

$$\text{read}(\lambda i. s, r) = s[i/r]$$

β -reduction

index terms	$t_I ::= \dots$
element terms	$t_E ::= \dots \mid \text{read}(t_A, t_I)$
array terms	$t_A ::= a \mid \text{write}(t_A, t_I, t_E) \mid \lambda i. t_E$

$$p = r \implies \text{read}(\text{write}(a, p, v), r) = v$$

$$\neg(p = r) \implies \text{read}(\text{write}(a, p, v), r) = \text{read}(a, r)$$

$$\text{read}(\lambda i. s, r) = s[i/r]$$

β -reduction

$\text{write}(a, p, v)$ could be **simulated** using $\lambda i. \text{ITE}(p = i, v, \text{read}(a, i))$

- Precisely model `memset` and `memcpy`

Uses of $\mathcal{T}_{\lambda A}$

- Precisely model `memset` and `memcpy`
- Summarize loops

Uses of $\mathcal{T}_{\lambda A}$

- Precisely model `memset` and `memcpy`
- Summarize loops
- Zero initialization of global variables

Uses of $\mathcal{T}_{\lambda A}$

- Precisely model `memset` and `memcpy`
- Summarize loops
- Zero initialization of global variables
- Zero initialization of fresh memory pages

- Precisely model `memset` and `memcpy`
- Summarize loops
- Zero initialization of global variables
- Zero initialization of fresh memory pages
- “Havoc” memory regions (volatile variables)

- Precisely model `memset` and `memcpy`
- Summarize loops
- Zero initialization of global variables
- Zero initialization of fresh memory pages
- “Havoc” memory regions (volatile variables)
- Model memory mapped I/O

- Precisely model `memset` and `memcpy`
- Summarize loops
- Zero initialization of global variables
- Zero initialization of fresh memory pages
- “Havoc” memory regions (volatile variables)
- Model memory mapped I/O
- Attaching metadata to memory regions (allocated, de-allocated, . . .)

Loop Summarization Using $\mathcal{T}_{\lambda A}$

- Broadly speaking:
 - loop iterations do not depend on **earlier** iterations
 - consecutive iterations update **consecutive** array locations

- Broadly speaking:
 - loop iterations do not depend on **earlier** iterations
 - consecutive iterations update **consecutive** array locations
- More precisely:
 - Induction variable i is **incremented by one** in each iteration
 - i^{th} iteration unconditionally **updates only** $a[i]$
 - **No other variable** declared outside the loop is **modified**
 - i^{th} iteration of the loop **may not** use elements of a that have been modified in earlier iterations

- Broadly speaking:
 - loop iterations do not depend on **earlier** iterations
 - consecutive iterations update **consecutive** array locations
- More precisely:
 - Induction variable i is **incremented by one** in each iteration
 - i^{th} iteration unconditionally **updates only** $a[i]$
 - **No other variable** declared outside the loop is **modified**
 - i^{th} iteration of the loop **may not** use elements of a that have been modified in earlier iterations
- Loops can often be **automatically transformed** into loops that satisfy these requirements

- Based on **reductions** to theories supported by SMT-solvers

- Based on **reductions** to theories supported by SMT-solvers
- One **quantifier-based** approach

- Based on **reductions** to theories supported by SMT-solvers
- One **quantifier-based** approach
- Two **quantifier-free** approaches
 - Eager reduction
 - Instantiation-based approach

Quantifier-Based Approach

- Replace $\lambda i. s$ by a fresh constant a_s

Quantifier-Based Approach

- Replace $\lambda i. s$ by a fresh constant a_s
- Add the constraint

$$\forall r. \text{read}(a_s, r) = s[i/r]$$

to the formula

Quantifier-Based Approach

- Replace $\lambda i. s$ by a **fresh constant** a_s
- **Add** the constraint

$$\forall r. \text{read}(a_s, r) = s[i/r]$$

to the formula

- Requires an SMT-solver that supports **quantifiers**

- Replace $\lambda i. s$ by a fresh constant a_s
- Add the constraint

$$\forall r. \text{read}(a_s, r) = s[i/r]$$

to the formula

- Requires an SMT-solver that supports quantifiers
- Does not provide a decision procedure in general

■ Replace `read(write(a, p, v), r)` by

$$\text{ITE}(p = r, v, \text{read}(a, r))$$

- Replace $\text{read}(\text{write}(a, p, v), r)$ by

$$\text{ITE}(p = r, v, \text{read}(a, r))$$

- Replace $\text{read}(\lambda i. s, r)$ by

$$s[i/r]$$

- Replace $\text{read}(\text{write}(a, p, v), r)$ by

$$\text{ITE}(p = r, v, \text{read}(a, r))$$

- Replace $\text{read}(\lambda i. s, r)$ by

$$s[i/r]$$

- $\mathcal{T}_{\lambda, A}$ axioms are applied eagerly

- Replace $\text{read}(\text{write}(a, p, v), r)$ by

$$\text{ITE}(p = r, v, \text{read}(a, r))$$

- Replace $\text{read}(\lambda i. s, r)$ by

$$s[i/r]$$

- $\mathcal{T}_{\lambda\mathcal{A}}$ axioms are **applied** eagerly
- Can be used in combination with **any** solver that supports $\mathcal{T}_{\mathcal{A}}$ and the index and element theories

Instantiation-Based Approach

- Replace $\lambda i. s$ by a fresh constant a_s

Instantiation-Based Approach

- Replace $\lambda i. s$ by a fresh constant a_s
- Add needed instantiations of

$$\forall r. \text{read}(a_s, r) = s[i/r]$$

to the formula

Instantiation-Based Approach

- Replace $\lambda i. s$ by a **fresh constant** a_s
- Add needed **instantiations** of

$$\forall r. \text{read}(a_s, r) = s[i/r]$$

to the formula

- Needed instantiations are **determined** by reads that “**depend**” on a_s

Instantiation-Based Approach

- Replace $\lambda i. s$ by a **fresh constant** a_s
- Add needed **instantiations** of

$$\forall r. \text{read}(a_s, r) = s[i/r]$$

to the formula

- Needed instantiations are **determined** by reads that “**depend**” on a_s
- Can be used in combination with **any** solver that supports \mathcal{T}_A and the index and element theories

- Done in the software bounded model checker [LLBMC](#)

Evaluation

- Done in the software bounded model checker [LLBMC](#)
- Uses [bitvectors](#) as index and element theories

- Done in the software bounded model checker [LLBMC](#)
- Uses [bitvectors](#) as index and element theories
- Applied on 81 benchmark programs
 - 67 programs produce λ -terms obtained from `memset` or `memcpy`
 - 14 program contain loops that can be summarized using λ -terms

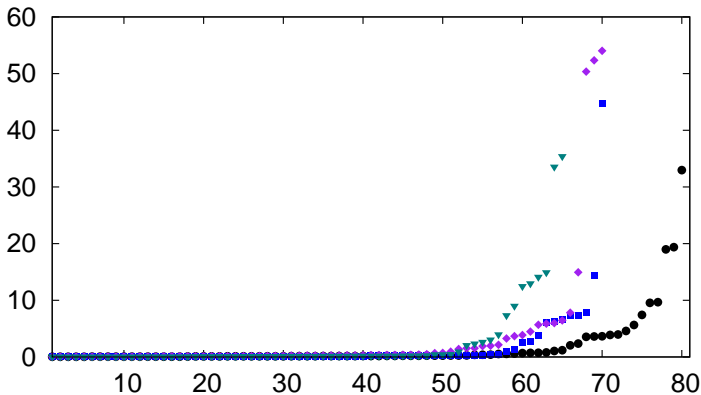
- Done in the software bounded model checker [LLBMC](#)
- Uses [bitvectors](#) as index and element theories
- Applied on 81 benchmark programs
 - 67 programs produce λ -terms obtained from `memset` or `memcpy`
 - 14 program contain loops that can be summarized using λ -terms
- Of the resulting formulas, 20 are satisfiable and 61 are unsatisfiable

- Done in the software bounded model checker `LLBMC`
- Uses `bitvectors` as index and element theories
- Applied on 81 benchmark programs
 - 67 programs produce λ -terms obtained from `memset` or `memcpy`
 - 14 program contain loops that can be summarized using λ -terms
- Of the resulting formulas, 20 are satisfiable and 61 are unsatisfiable
- Evaluated three `reductions` and `loop unrolling`
 - Quantifier-based approach using `Z3` and `CVC4`
 - Eager reduction and instantiation-based approach using `STP`, `Boolector`, `Z3`, and `CVC4`
 - Loop unrolling approach using `STP`, `Boolector`, `Z3`, and `CVC4`

SMT solver	Approach	Total Time	# Solved Formulas	# Timeouts	# Aborts
STP	Instantiation	206.034	80	1	–
STP	Eager	779.544	70	11	–
STP	Loops	670.526	70	6	5
Boolector	Instantiation	818.782	71	10	–
Boolector	Eager	986.751	70	11	–
Boolector	Loops	1139.483	61	15	5
Z3	Instantiation	948.365	67	13	1
Z3	Eager	1043.632	66	15	–
Z3	Quantifiers	1122.489	65	16	–
Z3	Loops	1619.583	53	23	5
CVC4	Instantiation	928.079	67	14	–
CVC4	Eager	1119.748	65	16	–
CVC4	Quantifiers	1407.118	54	21	6
CVC4	Loops	1552.698	56	19	6

Results

Instantiation (STP) ●
Eager (STP) ■
Loops (STP) ◆
Quantifiers (Z3) ▼



Conclusion and Future Work

- $\mathcal{T}_{\lambda\mathcal{A}}$ is a **useful, decidable extension** of $\mathcal{T}_{\mathcal{A}}$

Conclusion and Future Work

- $\mathcal{T}_{\lambda\mathcal{A}}$ is a **useful, decidable extension** of $\mathcal{T}_{\mathcal{A}}$
- Performs **better than unrolling** for
 - `memset` and `memcpy`
 - summarizable loops

- $\mathcal{T}_{\lambda\mathcal{A}}$ is a **useful, decidable extension** of $\mathcal{T}_{\mathcal{A}}$
- Performs **better than unrolling** for
 - memset and memcpy
 - summarizable loops
- Quantifier-free reductions perform **better than Z3's and CVC4's** reasoning involving **quantifiers**

- $\mathcal{T}_{\lambda\mathcal{A}}$ is a **useful, decidable extension** of $\mathcal{T}_{\mathcal{A}}$
- Performs **better than unrolling** for
 - memset and memcpy
 - summarizable loops
- Quantifier-free reductions perform **better than Z3's** and CVC4's reasoning involving **quantifiers**
- Integration into an SMT-solver using “**Lemmas-on-demand**”/“**lazy instantiation**” is the next step

`http://llbmc.org`