

Arrays Made Simpler : An Efficient, Scalable and Thorough Preprocessing

Logic for Programming, Artificial Intelligence and Reasoning

- | | | |
|--------------------------|---|-----------|
| Benjamin Farinier | — | CEA List |
| Robin David | — | Quarkslab |
| Sébastien Bardin | — | CEA List |
| Matthieu Lemerre | — | CEA List |

Software verification

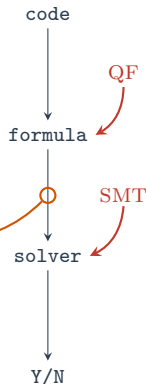
- More and more rely on decision procedures

Array theory

- Useful for modelling memory or data structures
- Bottleneck of resolution for large formulas (BMC, SE)

Our proposal

- FAS, an efficient **simplification** for array theory
- ⇒ Improve existing solvers



Two basic operations on arrays

- **Reading** in a at index $i \in \mathcal{I}$: $a[i]$
- **Writing** in a an element $e \in \mathcal{E}$ at index $i \in \mathcal{I}$: $a[i] \leftarrow e$

read : Array $\mathcal{I} \mathcal{E} \rightarrow \mathcal{I} \rightarrow \mathcal{E}$

write : Array $\mathcal{I} \mathcal{E} \rightarrow \mathcal{I} \rightarrow \mathcal{E} \rightarrow$ Array $\mathcal{I} \mathcal{E}$

$$\forall a i j e. (a[i] \leftarrow e)[j] = \begin{cases} e & \text{if } i = j \\ a[j] & \text{otherwise} \end{cases}$$

Prevalent in software analysis

- Modelling memory
- Abstracting data structure (map, queue, stack...)

Hard to solve

- NP-complete
- ROW may require case-splits

Unrolling-based verification techniques (BMC, SE)

- may produce huge formula
- high number of reads and writes

In some extremes cases, solvers may spend **hours** of these formulas

Without proper simplification,
array theory might become a bottleneck for resolution

What should we simplify ? **Read-Over-Write** (ROW) !

ROW simplification

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```



```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
assert (mem0 [1415] = 9265)
```

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```



```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
assert (mem0 [1415] = 9265)
```

These simplifications depend on two factors

- The equality check procedure
 \implies verify that $eax_0 = esp_0 - 16$
- The underlying representation of an array
 \implies remember that $mem_1 [esp_1 + 48] = 1415$

- ① Dedicated data structure
- ② Approximated equality check
- ③ Experimental evaluation
- ④ Conclusion

Section 1

Dedicated data structure

Standard ROW simplification (using list representation)

Array terms represented as write chains



How to update

Given a write of e at index i

- Add a new head node (i, e)

How to simplify ROW

Given a read at index j

- Let (i, e) be the head
- Is $i = j$ *valid*? If so return e
- Is $i \neq j$ *valid*? If so recurse
- Else stop

Standard ROW simplification (using list representation)

Array terms represented as write chains



How to update

Given a write of e at index i

- Add a new head node (i, e)

How to simplify ROW

Given a read at index j

- Let (i, e) be the head
- Is $i = j$ valid? If so return e
- Is $i \neq j$ valid? If so recurse
- Else stop

Good points

- Generic
- Full ROW simplification

Bad points

- n^2 worst case
- $k \cdot n$ if k -bounded, but then partial simplification
- Behave badly with WOW

Common preprocessing in SE when all indices are constants



Bad point

- Extremely specific

Good points

- $n \cdot \ln n$ worst case
- Full ROW simplification
- Behave well with WOW

So here we are

	list	k -list	map
generic	✓	✓	✗
ROW	✓	✓/✗	✓
WOW	✗	✗	✓
complexity	n^2	$k \cdot n$	$n \cdot \ln n$

We would like the best of both worlds!

	list	k -list	map
generic	✓	✓	✗
ROW	✓	✓/✗	✓
WOW	✗	✗	✓
complexity	n^2	$k \cdot n$	$n \cdot \ln n$

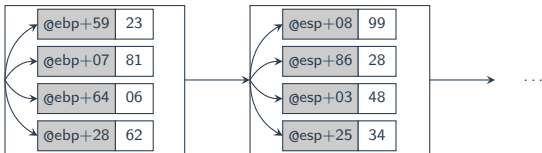
We would like the best of both worlds!

Our proposal : List-Map representation

- Sets of comparable indexes can be packed together
- Allows efficient search during ROW-simplification

⇒ Improve scalability

Improving scalability : list-map representation

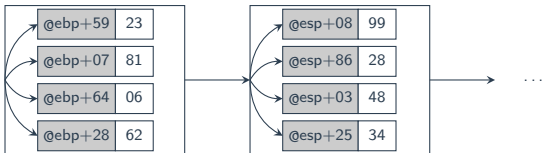


How to update

Given a write of e at index i

- Is i *comparable* with indices of elements in the head?
- If so add (i, e) in this map
- Else add a new head map containing only (i, e)

Improving scalability : list-map representation



How to update

Given a write of e at index i

- Is i *comparable* with indices of elements in the head?
- If so add (i, e) in this map
- Else add a new head map containing only (i, e)

How to simplify ROW

Given a read at index j

- Is j *comparable* with indices of elements in the head?
- If so, look for (i, e) with $i=j$
 - if succeed then return e
 - else recurse on next map
- Else stop

Improving scalability : list-map representation



How to update

Given a write of e at index i

- Is i *comparable* with indices of elements in the head ?
- If so add (i, e) in this map
- Else add a new head map containing only (i, e)

How to simplify ROW

Given a read at index j

- Is j *comparable* with indices of elements in the head ?
- If so, look for (i, e) with $i=j$
 - if succeed then return e
 - else recurse on next map
- Else stop

Section 2

Approximated equality check

Propagate “variable+constant” terms

- If $y \triangleq z + 1$ then $x \triangleq y + 2 \rightsquigarrow x \triangleq z + 3$
- Together with associativity, commutativity...

\implies Reduce the number of bases

Propagate “variable+constant” terms

- If $y \triangleq z + 1$ then $x \triangleq y + 2 \rightsquigarrow x \triangleq z + 3$
- Together with associativity, commutativity...

\implies Reduce the number of bases

if $u \triangleq v$ then	$u + k$	\rightsquigarrow	$v + k$	alias inlining
if $u \triangleq v + l$ then	$u + k$	\rightsquigarrow	$v + (k + l)$	base/offset inlining
	$-(x + k)$	\rightsquigarrow	$(-k) - x$	constant negation
	$(x + k) + l$	\rightsquigarrow	$x + (k + l)$	constant packing
	$(x + k) + y$	\rightsquigarrow	$(x + y) + k$	constant lifting
	$(x + k) + (y + l)$	\rightsquigarrow	$(x + y) + (k + l)$	base/offset addition
	$(x + k) - (y + l)$	\rightsquigarrow	$(x - y) + (k - l)$	base/offset subtraction
		\dots		

Propagate “variable+constant” terms

- If $y \triangleq z + 1$ then $x \triangleq y + 2 \rightsquigarrow x \triangleq z + 3$
- Together with associativity, commutativity...

⇒ Reduce the number of bases

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```



```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```

Propagate “variable+constant” terms

- If $y \triangleq z + 1$ then $x \triangleq y + 2 \rightsquigarrow x \triangleq z + 3$
- Together with associativity, commutativity...

⇒ Reduce the number of bases

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```



```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp0 - 64 + 48]  
assert (mem1 [eax0] = 9265)
```

Propagate “variable+constant” terms

- If $y \triangleq z + 1$ then $x \triangleq y + 2 \rightsquigarrow x \triangleq z + 3$
- Together with associativity, commutativity...

\implies Reduce the number of bases

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```



```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
eax0  $\triangleq$  mem1 [esp0 - 16]  
assert (mem1 [eax0] = 9265)
```

Propagate “variable+constant” terms

- If $y \triangleq z + 1$ then $x \triangleq y + 2 \rightsquigarrow x \triangleq z + 3$
- Together with associativity, commutativity...

\implies Reduce the number of bases

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```



```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
eax0  $\triangleq$  1415  
assert (mem1 [eax0] = 9265)
```

Propagate “variable+constant” terms

- If $y \triangleq z + 1$ then $x \triangleq y + 2 \rightsquigarrow x \triangleq z + 3$
- Together with associativity, commutativity...

\implies Reduce the number of bases

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```



```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
eax0  $\triangleq$  1415  
assert (mem1 [1415] = 9265)
```


Propagate “variable+constant” terms

- If $y \triangleq z + 1$ then $x \triangleq y + 2 \rightsquigarrow x \triangleq z + 3$
- Together with associativity, commutativity...

\implies Reduce the number of bases

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```



```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
assert (mem1 [1415] = 9265)
```

How to get ride of this declaration ?

Associate to every indices i an abstract domain i^\sharp

- If $i^\sharp \sqcap j^\sharp = \perp$ then $(a[i] \leftarrow e)[j] = a[j]$
- Integrated in the list-map representation

\implies Prove disequality between different bases

Associate to every indices i an abstract domain i^\sharp

- If $i^\sharp \sqcap j^\sharp = \perp$ then $(a[i] \leftarrow e)[j] = a[j]$
- Integrated in the list-map representation

\implies Prove disequality between different bases

$$c^\sharp = [c, c]$$

for any constant c

$$v^\sharp = [m_i, M_j]$$

if $i \leq v \leq j$

$$\begin{aligned} (\text{extract}_{l,h} i)^\sharp &= [0, 2^{h-l+1} - 1] \\ &= [\text{extract}_{l,h}(m_i), \text{extract}_{l,h}(M_j)] \\ &= [0, \text{extract}_{l,h}(M_j)] \\ &\sqcup [\text{extract}_{l,h}(m_i), 2^{h-l+1} - 1] \end{aligned}$$

if $(M_j \gg l) - (m_i \gg l) \geq 2^{h-l+1}$
if $\text{extract}_{l,h}(M_j) \geq \text{extract}_{l,h}(m_i)$
otherwise

$$\begin{aligned} (i + j)^\sharp &= [m_i + m_j, M_i + M_j] \\ &= [m_i + m_j - 2^N, M_i + M_j - 2^N] \\ &= [m_i + m_j - 2^N, 2^N - 1] \\ &\sqcup [0, M_i + M_j - 2^N] \end{aligned}$$

if $M_i + M_j < 2^N$
if $m_i + m_j \geq 2^N$
otherwise

Associate to every indices i an abstract domain $i^\#$

- If $i^\# \sqcap j^\# = \perp$ then $(a[i] \leftarrow e)[j] = a[j]$
- Integrated in the list-map representation

\implies Prove disequality between different bases

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```



```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
assert (mem1 [1415] = 9265)  
esp0# = [61441, 615535]
```

Associate to every indices i an abstract domain $i^\#$

- If $i^\# \sqcap j^\# = \perp$ then $(a[i] \leftarrow e)[j] = a[j]$
- Integrated in the list-map representation

\implies Prove disequality between different bases

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```



```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
assert (mem1 [1415] = 9265)  
  
(esp0 - 16)# = [61425, 615519]
```

Associate to every indices i an abstract domain $i^\#$

- If $i^\# \sqcap j^\# = \perp$ then $(a[i] \leftarrow e)[j] = a[j]$
- Integrated in the list-map representation

\implies Prove disequality between different bases

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```



```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
assert (mem1 [1415] = 9265)
```

$1415^\# = [1415, 1415]$

Associate to every indices i an abstract domain $i^\#$

- If $i^\# \sqcap j^\# = \perp$ then $(a[i] \leftarrow e)[j] = a[j]$
- Integrated in the list-map representation

\implies Prove disequality between different bases

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```



```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
assert (mem0 [1415] = 9265)  
  
(esp0 - 16)#  $\sqcap$  1415# =  $\perp$ 
```

Associate to every indices i an abstract domain $i^\#$

- If $i^\# \sqcap j^\# = \perp$ then $(a[i] \leftarrow e)[j] = a[j]$
- Integrated in the list-map representation

\implies Prove disequality between different bases

```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
mem1  $\triangleq$  mem0 [esp0 - 16]  $\leftarrow$  1415  
esp1  $\triangleq$  esp0 - 64  
eax0  $\triangleq$  mem1 [esp1 + 48]  
assert (mem1 [eax0] = 9265)
```



```
esp0 : BitVec16  
mem0 : Array BitVec16 BitVec16
```

```
assert (esp0 > 61440)  
assert (mem0 [1415] = 9265)
```


Section 3

Experimental evaluation

FAS implemented as

- preprocessor for SMT formulas
- QF_ABV theory

The implementation comprises

- 6,300 lines of OCaml
- integrated into the TFML formula preprocessing engine
- part of the BINSEC Symbolic Execution (SE) tool

Experiments are carried out on

- Intel(R) Xeon(R) CPU E5-2660 v3 @ 2,60GHz
- three of the best SMT solvers for the QF_ABV theory
Boolector, Yices and Z3

Impact of the simplification : medium-size formulas

- 6,590 × 3 medium-size formulas from static SE
- TIMEOUT = 1,000 seconds

		simpl. time	#TIMEOUT and resolution time						#ROW
			Boolector		Yices		Z3		
concrete	default	61	0	163	2	69	0	872	866,155
	FAS	85	0	94	2	68	0	244	1,318
	FAS-itv	111	0	94	2	68	0	224	1,318
interval	default	65	0	2,584	2	465	31	155,992	866,155
	FAS	99	0	2,245	2	487	25	126,806	531,654
	FAS-itv	118	0	755	2	140	14	37,269	205,733
symbolic	default	61	0	6,173	3	1,961	65	305,619	866,155
	FAS	91	0	6,117	3	1,965	66	158,635	531,654
	FAS-itv	111	0	4,767	2	1,108	43	80,569	295,333

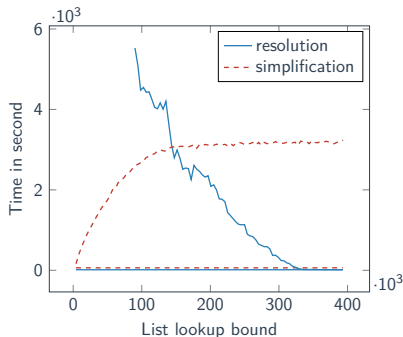
Impact of the simplification : very large formulas

- 29 × 3 very large formulas from dynamic SE
- TIMEOUT = 1,000 seconds

		simpl. time	#TIMEOUT and resolution time						#ROW
			Boolector		Yices		Z3		
concrete	default	44	10	159	4	1,098	26	3.33	1,120,798
	FAS-list	1,108	8	845	4	198	10	918	456,915
	FAS	196	8	820	4	196	10	922	456,915
	FAS-itv	210	4	654	1	12	4	1,120	0
interval	default	44	12	131	12	596	27	0.19	1,120,798
	FAS-list	222	12	129	12	595	26	236	657,594
	FAS	231	12	129	12	597	26	291	657,594
	FAS-itv	237	12	58	12	28	19	81	651,449
symbolic	default	40	12	1,522	12	1,961	27	0.13	1,120,798
	FAS-list	187	11	1,199	12	2,018	26	486	657,594
	FAS	194	11	1,212	12	2,081	26	481	657,594
	FAS-itv	200	11	1,205	12	2,063	26	416	657,594

Focus on specific case : the ASPack example

- Huge formula obtained from the ASPack packing tool
- 293 000 ROWs
- 24 hours of resolution !



Using FAS

- #ROW reduced to 2 467
- 14 sec for resolution
- 61 sec for preprocessing

Using list representation

- Same result with a bound of 385 024 and beyond...
- ...but 53 min preprocessing

Section 4

Conclusion

Software verification

- More and more rely on decision procedures

Array theory

- Useful for modelling memory or data structures
- Bottleneck of resolution for large formulas (BMC, SE)

Our proposal

- Efficient **simplification** for array theory

⇒ Improve existing solvers

Future work

- Deeper integration inside a dedicated array solver
- Adding more expressive domain reasoning