

Fine-Grained Formal Specification and Analysis of Buddy Memory Allocation in Zephyr RTOS

北京航空航天大学计算机新技术研究所
The Institute of Advanced Computing Technology



Authors: **Feng Zhang¹, Yongwang Zhao^{*,1}, Dianfu Ma¹, Wensheng Niu²**

*. Corresponding Author

1. School of Computer Science and Engineering, Beihang University, China

2. Aeronautical Computing Technique Research Institute, Xi'an, China

Contents

- 1. Introduction
- 2. Buddy Memory Allocation Algorithm in Zephyr
- 3. Fine-Grained Formal Specification in Isabell/HOL
- 4. Formal Proof
- 5. Results and Discussions
- 6. Conclusions

1. Introduction - Abstract

- Memory management (MM) is a critical component of OS
- **Bugs** in MM may crash OS or the whole critical system
- This paper presents a case study of formal verification on the **buddy** memory allocation component of the Zephyr RTOS:
 - Provide **Fine-Grained** formal specification in Isabelle/HOL
 - Conduct Formal proof using the **interactive** prover in Isabelle
 - Find **two flaws** in the C code when executing sequentially

1. Introduction – Research Status

- Verification of the **TLSF** algorithm in Event-B:
 - Only verifies an abstract specification at the requirement level
 - not check consistency between elements in the data structure
- **seL4** pushes the memory allocation outside of the kernel
- Yu et al. introduce a low-level language **CAP** (certified assembly programming) in Coq
 - build certified programs
 - present a certified library for dynamic storage allocation
 - not a kernel's component but a certified memory library
 - 75 lines C code

1. Introduction – summary

- We create a fine-grained formal specification:
 - All the elements of the data structure
 - All the operations (initialization, allocation and release)
 - System clocks and simple kernel scheduling
 - The execution of memory allocation is preemptive
- We concentrate in five types of critical properties:
 - Invariants
 - Correctness of doubly linked lists
 - Functional correctness of events
 - Conformity of event specifications to kernel requirements
 - Livelock-free of the system specification.

Contents

- 1. Introduction
- 2. Buddy Memory Allocation Algorithm in Zephyr
- 3. Fine-Grained Formal Specification in Isabell/HOL
- 4. Formal Proof
- 5. Results and Discussions
- 6. Conclusions



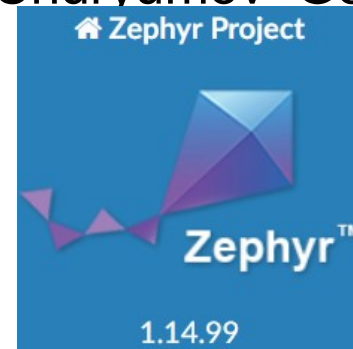
2.1 – Zephyr Project

- Zephyr Project is a **Linux** Foundation Project
- Be perfect for building simple connected **sensors**:
 - up to modems and small IoT wireless gateways
 - Built with **safety** and **security** in mind
 - Cross-architecture with growing developer tool support
 - Complete, fully integrated, highly configurable, **modular** for flexibility, better than roll-your-own
 - Product development ready
 - Permissively licensed



2.2 – Zephyr OS Kernel

- Derived from **Wind River**'s commercial Microkernel Profile
- Microkernel Profile has evolved over 20 years from DSP RTOS technology known as Virtuoso
- Used in several commercial applications:
 - satellites, military command and control communications, radar, telecommunications and image processing
 - successful **Philae** Landing on Comet Churyumov–Gerasimenko and the accompanying Rosetta Orbiter



2.2 – Buddy Memory Allocation Algorithm in Zephyr Kernel

➤ (1) Pool and block Initialization

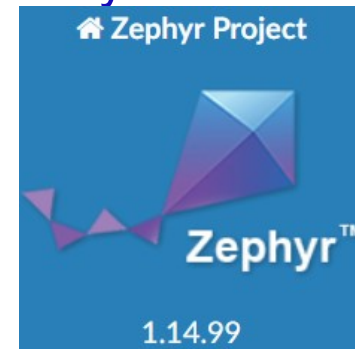
- only be defined and initialized at compile time

➤ (2) Block Allocation

- Quad-Partitioning: iteratively partitioning larger blocks into smaller quad-ones

➤ (3) Block Release

- Immediately, automatically, and recursively combining smaller blocks into bigger ones



2.2 – Buddy Memory Allocation Algorithm in Zephyr Kernel

Algorithm 1 pool_alloc (p, size, block)

```

Input: p: the requester;
        block: information
Output: ZERO: success;
        1: alloc_l = -1; free_l = 0;
        2: for i=0; i < p.n_level; i++
        3:   if block size at level i is greater than or equal to size
        4:     break
        5:   end if
        6:   alloc_l = i
        7:   if ! level_empty(i)
        8:     free_l = i
        9:   end if
        10: end for
        11: if alloc_l < 0 || free_l == 0
        12:   return NOMEM
        13: end if
        14: blk ← address of block at level alloc_l and free_l
        15: if ! blk then
        16:   return EAGAIN
        17: end if
        18: for from_l = free_l; from_l < alloc_l; from_l++
        19:   partitioning block at level from_l
        20:   blk ← address of block at level from_l
        21: end for
        22: block ← (blk, p[from_l][blk])
        23: return 0
    
```

Algorithm 2 k_mem_pool_alloc (p, size, block, timeout)

```

Input: p: the requester;
        block: information;
        timeout: no_wait
Output: ZERO: success;
        1: if size > p.max_block_size
        2:   return NOMEM
        3: else
        4:   if timeout > 0
        5:     end = cur_time + timeout
        6:   end if
        7:   while (1) do
        8:     ret = pool_alloc(p, size, block)
        9:     if ret == 0
        10:      return ret
        11:     end if
        12:     SCHEDULE
        13:     if timeout > 0
        14:       timeo = end - cur_time
        15:       if timeo <= 0 then
        16:         break
        17:       end if
        18:     end if
        19:   end while
        20:   return TIMEOUT
        21: end if
    
```

Algorithm 3 free_block (p, level, lsize, bn)

```

Input: p: the released pool; level: the released level number;
        lsize: an array contains each level's block size; bn: the released block number
Output: void
        1: set_bit(p, level, bn)
        2: if level > 0 && partner of bn are all free then
        3:   for i=0; i < 4; i++ do
        4:     clear_bit(p, level, bn+i)
        5:     if (bn & ~3) + i ≠ bn then
        6:       remove block ((bn & ~3) + i) from its free list
        7:     end if
        8:   end for
        9:   free_block(p, level-1, lsize, bn/4)
        10:  return
        11: end if
        12: append the the released block (p, level, bn) into its free list
    
```

Contents

- 1. Introduction
- 2. Buddy Memory Allocation Algorithm
- 3. Fine-Grained Formal Specification in Isabell/HOL
- 4. Formal Proof
- 5. Results and Discussions
- 6. Conclusions



3 – Fine-Grained Formal Specification

➤ A. State Machine

- The state is defined as a **record** $StateD$
- the initial state s_0
- state-transition functions φ

Definition 1: State machine of the buddy memory allocation Component $\mathcal{M} = \langle \mathcal{S}, \mathcal{E}, \varphi, s_0 \rangle$ is a tuple, where \mathcal{S} is the state space, \mathcal{E} is a set of event labels, $s_0 \in \mathcal{S}$ is the initial state, and $\varphi: \mathcal{E} \rightarrow \mathbb{P}(\mathcal{S} \times \mathcal{S})$ is a set of state-transition functions.



3 – Fine-Grained Formal Specification

➤ B. Data Structure

```

struct k_mem_block_id {
    u32_t pool : 8;
    u32_t level : 4;
    u32_t block : 20; };

struct k_mem_block {
    void *data;
    struct k_mem_block_id id; };

struct k_mem_pool_lvl {
    union {
        u32_t *bits_p;
        u32_t bits;
    };
    sys_dlist_t free_list; };

struct k_mem_pool {
    void *buf;
    size_t max_sz;
    u16_t n_max;
    u8_t n_levels;
    u8_t max_inline_level;
    struct k_mem_pool_lvl *levels;
    _wait_q_t wait_q; };
    
```

```

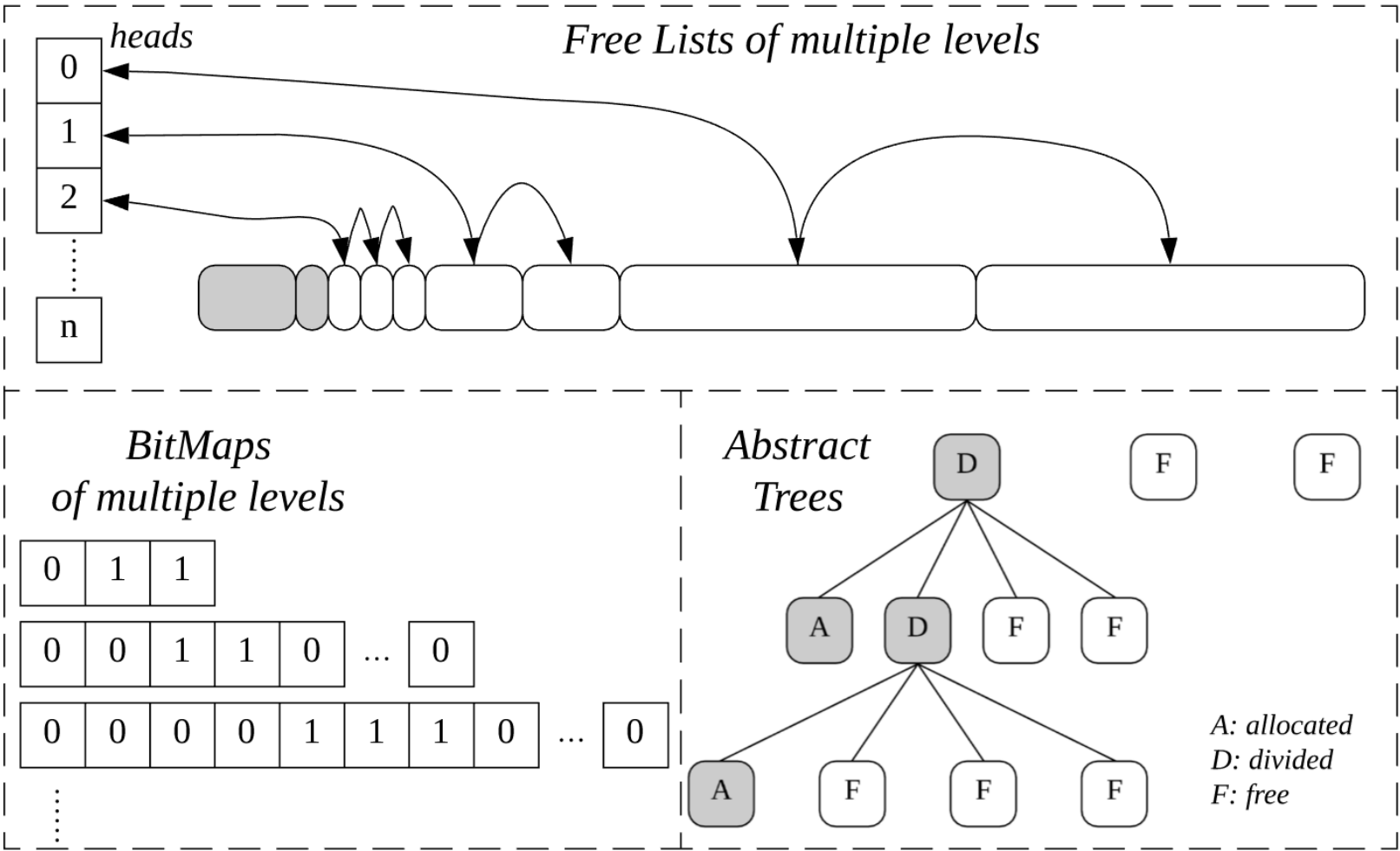
type_synonym struct_block_id = "pool_num × level_num × block_num"
type_synonym struct_block = "struct_block_id × addr"
type_synonym BlockD = "(struct_block_id × block_state_type) tree"
type_synonym struct_pool_lvl = "bitMap × ref_freelist"
record PoolD = l0blist :: "BlockD list"
    name :: string
    max_sz :: nat
    nmax :: nat
    n_levels :: nat
    max_inline_level :: max_inline_lsz
    levels :: "struct_pool_lvl list"

record heap = refs :: "ref set"
    addrS :: "addr set"
    addr2bid :: "addr → struct_block_id"
    ref2bid :: "sys_dnode_t → struct_block_id"
    head_next :: "ref → ref"
    tail_prev :: "ref → ref"

record StateD = globals :: heap
    poolsD :: "PoolD list"
    curD :: "Thread option"
    irq :: bool
    tickD :: nat
    t_stateD :: "Thread ⇒ thread_state_type"
    waitqD :: "Thread → pool_num"
    alloc_context ::
        "Thread → (pool_num × request_size × timeout × end_time)"
    
```

3 – Fine-Grained Formal Specification

➤ B. Data Structure



3 – Fine-Grained Formal Specification

➤ C. Event Specification

□ system behaviors based on Zephyr characteristics

- system clocks *time_tick*
- the thread scheduling *schedule*

□ actions operated on memory pools and blocks

- pool and block initializations
- block allocations
- block release

3 – Fine-Grained Formal Specification

➤ C. Event Specification

□ system behavior

- system clocks
- the thread scheduler

□ actions operation

- pool and block allocation
- block allocation
- block release

Algorithm 2 `k_mem_pool_alloc` ($p, size, block, timeout$)

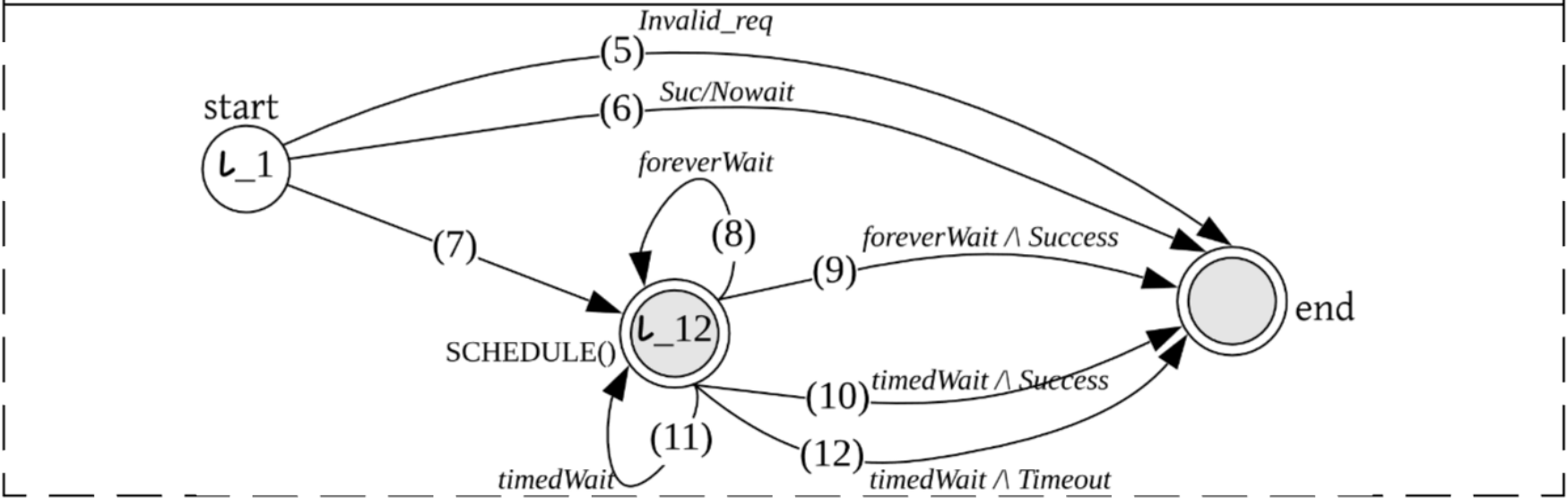
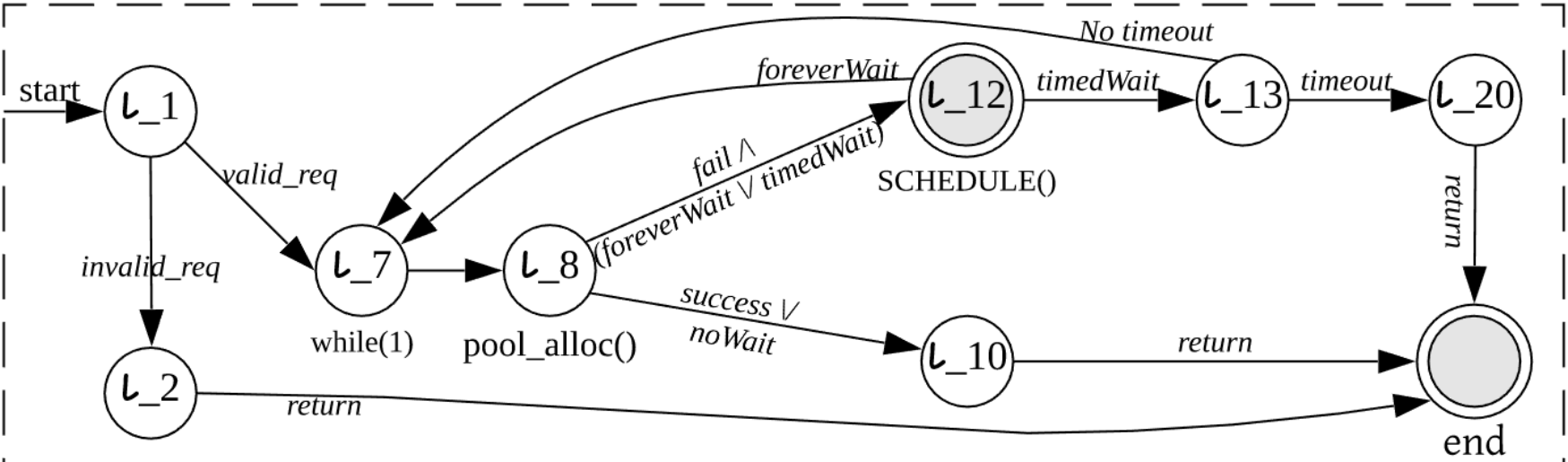
Input: p : the requested pool; $size$: the requested size
 $block$: information of the allocated block;
 $timeout$: `no_wait`, `forever` or `timed_wait`

Output: `ZERO`: successful allocation; `NOMEM`: no memory; `TIMEOUT`: timeout

```
1: if  $size > p.max\_sz$  then
2:   return NOMEM
3: else
4:   if  $timeout > 0$  then
5:      $end = current\_time + timeout$ 
6:   end if
7:   while (1) do
8:      $ret = pool\_alloc(p, size, block)$ 
9:     if  $ret == 0 \parallel timeout == no\_wait$  then
10:      return  $ret$ 
11:    end if
12:    SCHEDULE();
13:    if  $timeout \neq forever$  then
14:       $timeout = end - current\_time$ 
15:      if  $timeout \leq 0$  then
16:        break
17:      end if
18:    end if
19:  end while
20:  return TIMEOUT
21: end if
```


3 – Fine-Grained Formal Specification

➤ C. Event Specification



3 – Fine-Grained Formal Specification

➤ D. State Space

type_synonym Trace = “StateD list”

inductive_set TraceSpace :: “Trace set”

definition “reachable s t =

$(\exists ts \in \text{TraceSpace}. ts \neq [] \wedge \text{last } ts = t)$ ”

definition “ReachStates $\equiv \{s. \text{reachable } s_0 s\}$ ”

Contents

- 1. Introduction
- 2. Buddy Memory Allocation Algorithm in Zephyr
- 3. Fine-Grained Formal Specification in Isabell/HOL
- 4. Formal Proof
- 5. Results and Discussions
- 6. Conclusions



4 – Formal Proof

➤ 4.1 Invariants - Consistency of Data Structure

- ***bitMap_freelistS s*** specifies the consistency between bit_maps and free lists
- ***bitMap_treeS s*** specifies the consistency between bit_maps and abstract trees.

“**definition** $\text{Inv_Bitmap_freelist_tree } s \equiv \text{bitMap_freelistS } s \wedge \text{bitMap_treeS } s$ ”

Table 3 The Specification *bitMap_freelistS*

definition “ $\text{bitMap_freelist } pn \ p \ h \equiv \forall ln < n_levels \ p. \forall bn < (nmax \ p) * (4^{\wedge} ln). \text{let } \text{freel_ln} = \text{snd}((\text{levels } p)!ln); \text{ndref} = (\text{case } (\text{get_refBYbid } (\text{ref2bid } h) \ (pn,ln,bn)) \ \text{of } \text{Some } \text{ndr} \Rightarrow \text{ndr}); \text{bmGet} = \text{get_bit_ptr } (\text{fst}((\text{levels } p)!ln)) \ bn \ \text{in } \text{fst } \text{bmGet} \ ! \ \text{snd } \text{bmGet} \ \longleftrightarrow \ \text{nodeINfreelist } (\text{head_next } h) \ \text{freel_ln } \ \text{ndref} \neq \text{None}$ ”

definition “ $\text{bitMap_freelistS } s \equiv \forall pn < \text{length}(\text{poolsD } s). \text{let } p = (\text{poolsD } s)!pn; h = \text{globals } s \ \text{in } \text{bitMap_freelist } pn \ p \ h$ ”

4 – Formal Proof

➤ 4.2 Correctness of Doubly Linked Lists

- The pointer in C is specified as a *ref* in Isabelle
- $ref = (UNIV::nat\ set)$
- `head_next` :: “*ref* => *ref*”
- `tail_prev` :: “*ref* => *ref*”

```
struct _dnode {
  union {
    struct _dnode *head; /* ptr to head of list (sys_dlist_t) */
    struct _dnode *next; /* ptr to next node (sys_dnode_t) */
  };
  union {
    struct _dnode *tail; /* ptr to tail of list (sys_dlist_t) */
    struct _dnode *prev; /* ptr to previous node (sys_dnode_t) */
  };
};
typedef struct _dnode sys_dlist_t;
typedef struct _dnode sys_dnode_t;
```

4 – Formal Proof

➤ 4.2 Correctness of Doubly Linked Lists

- Length of a dilist
- Validity of a node
- Validity of a dlist
- Validity of appending actions

4 – Formal Proof

➤ 4.3 Functional Correctness of Events

- $\{P\} C \{Q\}$
- Our specifications are all total correctness specifications
- terminations are ensured by using the *primrec*, *fun*, *function* and *definition*

Lemma 11. *correctness of function allocL_freeL*

$$\begin{aligned} & \{ | \text{lsz} = n_levels \ p \wedge n_levels \ p > 0 \wedge \text{max_sz} \ p > 0 \}; \\ & \text{freel} = \text{snd}((\text{levels } p)!(\text{nat } (\text{snd } \text{alfl}))) \ | \} \\ \text{alfl} & = \mathbf{allocL_freeL} \ p \ h \ \text{rsz} \ (-1, -1) \ \text{lsz} \\ & \{ | \text{fst } \text{alfl} > -1 \wedge \text{snd } \text{alfl} > -1 \longrightarrow \\ & (\text{fst } \text{alfl} \geq \text{snd } \text{alfl} \wedge \text{fst } \text{alfl} \geq 0 \wedge \text{snd } \text{alfl} \geq 0 \wedge \\ & \text{nat}(\text{fst } \text{alfl}) < n_levels \ p \wedge \text{nat}(\text{snd } \text{alfl}) < n_levels \ p \wedge \\ & \neg \text{level_empty } (\text{head_next } h) \ \text{freel}) \ | \} \end{aligned}$$

4 – Formal Proof

➤ 4.4 Conformity of Event Specifications to Kernel Requirements

Lemma ExFreeNdPartners:

```
" $\forall s$  pn bn. let pn_s = (poolsD s)!pn; ln_i = n_levels pn_s - i; lvls_s = levels pn_s in
poolsD s  $\neq$  []  $\wedge$  pn < length(poolsD s)  $\wedge$  0 < i  $\wedge$  i  $\leq$  n_levels pn_s  $\wedge$  bn < (nmax pn_s) * (4^ln_i)  $\wedge$ 
reachableD s0D s  $\wedge$  fst (lookup_tree ((l0blist pn_s)!(bn div (4^ln_i))) (pn,ln_i,bn))  $\longrightarrow$ 
(if  $\neg$ partner_bitsLn lvls_s ln_i bn then
   $\exists t$ . let pn_t = (poolsD t)!pn in reachableD s t  $\wedge$ 
    ( $\forall$ ln' < n_levels pn_t.  $\forall$ bn' < (nmax pn_t) * (4^ln')).
    let blk_n_t' = lookup_tree ((l0blist pn_t)!(bn' div (4^ln'))) (pn,ln',bn');
      blk_n_s' = lookup_tree ((l0blist pn_s)!(bn' div (4^ln'))) (pn,ln',bn') in
      if ln'=ln_i  $\wedge$  bn'=bn then fst blk_n_t'  $\wedge$  isLeaf(snd blk_n_t')  $\wedge$  bstate (snd blk_n_t')=FREE
      else blk_n_t'=blk_n_s')
else  $\exists t$  j. let ln_j = ln_i - j;
  pn_t = (poolsD t)!pn in reachableD s t  $\wedge$  0 < j  $\wedge$  j  $\leq$  ln_i  $\wedge$ 
  ( $\forall$ ln' < n_levels pn_t.  $\forall$ bn' < (nmax pn_t) * (4^ln')).
  let blk_n_t' = lookup_tree ((l0blist pn_t)!(bn' div (4^ln'))) (pn,ln',bn');
    blk_n_s' = lookup_tree ((l0blist pn_s)!(bn' div (4^ln'))) (pn,ln',bn');
    bn_ln'_l = (bn div (4^j)) * (4^(ln'-ln_j));
    bn_ln'_r = (bn div (4^j) + 1) * (4^(ln'-ln_j)) in
    if ln'=ln_j  $\wedge$  bn'=bn div (4^j) then
      fst blk_n_t'  $\wedge$  isLeaf(snd blk_n_t')  $\wedge$  bstate (snd blk_n_t')=FREE
    else if ln' > ln_j  $\wedge$  (bn'  $\geq$  bn_ln'_l  $\wedge$  bn' < bn_ln'_r) then  $\neg$  fst blk_n_t'
    else blk_n_t'=blk_n_s'))"
```


Contents

- 1. Introduction
- 2. Buddy Memory Allocation Algorithm in Zephyr
- 3. Fine-Grained Formal Specification in Isabell/HOL
- 4. Formal Proof
- 5. Results and Discussions
- 6. Conclusions

5 –Results and Discussions

➤ A. Evaluation

- 600 lines C
- 800 lines specification:
 - 109 functions/definitions
 - 12 primary events
- 9400 lines proof: 338 lemmas

C code	Isabelle Code				total lines in Isabelle
	Specification		Proof		
.c	function/ definition	LOC	theorem/ lemma	LOP	
600	109	800	338	9400	10200

5 –Results and Discussions

- B. Results of formal analysis: fine two flaws
 - ❑ Return code not conform to the kernel requirement
 - ❑ Application thread will fall into live lock.



5

➤ B. Results of

□ Return code

□ Application

```
290: int k_mem_pool_alloc(struct k_mem_pool *p, struct k_mem_block *block,
291:                      size_t size, s32_t timeout)
292: {
293:     int ret, key;
294:     s64_t end = 0;
295:
296:     __ASSERT(!_is_in_isr() && timeout != K_NO_WAIT, "");
297:
298:     if (timeout > 0) {
299:         end = _tick_get() + _ms_to_ticks(timeout);
300:     }
301:
302:     while (1) {
303:         ret = pool_alloc(p, block, size);
304:
305:         if (ret == 0 || timeout == K_NO_WAIT ||
306:             ret == -EAGAIN || (ret && ret != -ENOMEM)) {
307:             return ret;
308:         }
309:
310:         key = irq_lock();
311:         _pend_current_thread(&p->wait_q, timeout);
312:         _Swap(key);
313:
314:         if (timeout != K_FOREVER) {
315:             timeout = end - _tick_get();
316:
317:             if (timeout < 0) {
318:                 break;
319:             }
320:         }
321:     } /* end while 1 */
322:     return -EAGAIN;
323: } /* end k_mem_pool_alloc */
324:
325:
```

Contents

- 1. Introduction
- 2. Buddy Memory Allocation Algorithm in Zephyr
- 3. Fine-Grained Formal Specification in Isabell/HOL
- 4. Formal Proof
- 5. Results and Discussions
- 6. Conclusions

6 – Conclusions

- We will perform formal analysis on the **concurrent characteristics** of the OS kernel
- For about 600lines C, our work consists of about 10200 lines of Isabelle
- Find two flaws in C code when executing sequentially



Thank you