



Experimental dependability evaluation of memory manager in the real-time operating system

Pawel Pisarczyk

Institute of Computer Science, Warsaw University of Technology, Warsaw 00-605 Nowowiejska 15/19, Poland
P.Pisarczyk@ii.pw.edu.pl

Abstract. The paper presents results of experimental dependability evaluation of the Phoenix-RTOS operating system. Experiments are conducted using a self-developed testing environment and a kernel fault injector. Dependability evaluation is the last stage of a system development process. Results will be used in the future research to propose the dependable memory manager.

1 Introduction

Dependability evaluation of operating systems is the focus of many research jobs. Developed methodologies are based mainly on software-fault injection technique. In papers [2, 8] new methods of dependability evaluation for microkernel-based operating systems are proposed. In paper [5] experimental environment for operating system with monolithic kernel and results of experimental dependability evaluation for Linux kernel are presented. The paper [6] presents model of error propagation between device drivers and operating system services. New measures characterizing error propagation (exposure and diffusion) have been introduced. Results for experimental estimation of introduced measures have been presented. In parallel to performing the presented works focused on dependability evaluation of popular operating systems, new architectures of mission critical operating systems have been proposed. The main goal of these architectures is to maximize efficiency of error detection and to eliminate fail silence violation. In paper [1, 3] the time-triggered architecture and results of experimental dependability evaluation for the architecture implementation (MARS operating system) have been presented. Architecture of MARS operating system based on hardware redundancy and dedicated communication interface. Time triggered architecture significantly decreases number of system applications.

Presented works, related to dependability evaluation of popular operating systems, are focused only on selected modules and system services provided for applications. Proposed operating system models do not include many functional dependencies between subsystems and error propagation is evaluated using simple metrics. Presented results have no influence on development of new methods and algorithms for error detection and propagation barriers embedded into operating systems. New algorithms for kernel-level error detection and fault tolerance have been only proposed for MARS operating system.

Phoenix-RTOS operating system developed by the author is an open-source real-time operating system for commonly used event-driven applications. System provides all mechanisms (like resource protection, threads and processes, IPC, IP networking and graphical interface) required by real industrial applications. Development process of every subsystem has been extended with dependability evaluation experiments. Results of these experiments allow to propose new error detection and fault tolerance algorithms integrated with the operating system kernel. Detailed error propagation analysis will be used to propose propagation barriers.

This paper presents results of experimental dependability evaluation of the memory management subsystem of Phoenix-RTOS real-time kernel Section 2 presents briefly Phoenix-RTOS operating system and memory management. Section 3 shows the used error model and introduced measures characterizing the system behavior after the fault injection. Section 4 describes experimental setup and developed fault injector embedded into the operating system kernel. Section 5 presents the analysis of results of conducted experiments. Section 6 presents conclusions.

2 Phoenix-RTOS

Phoenix-RTOS operating system has been developed in the Institute of Computer Science in Warsaw University of Technology in years 2005-2006 as the successor of the prototype operating system Phoenix [7]. System has been developed as a result of work which goal was to prepare the prototype for real industrial applications. Phoenix-RTOS is the time sharing real-time operating system intended for embedded

systems. System executes threads located in separated address spaces associated with processes. The system has been implemented for IA32 based PC computers used in embedded applications. Port for ARM7 architecture is available. Currently system is ported to PowerPC architecture. Phoenix-RTOS has been developed for real industrial needs which is the health monitoring appliance of one of a Polish medical equipment vendor. This paper is not focused on Phoenix-RTOS architecture therefore only the memory subsystem will be briefly presented.

Virtual memory management subsystem consists of five parts: hardware abstraction layer V_{hal} , page allocator V_{pg} , kernel virtual space allocator V_{kmap} , kernel heap allocator V_{kmal} and process virtual space allocator V_{map} . Memory manager functions are presented in the table 1.

Table 1. Virtual memory subsystem functions

Function	Byte size	Instructions	Part
<code>pmap_enter</code>	610	185	V_{hal}
<code>vm_pageAlloc</code>	953	291	V_{pg}
<code>vm_pageFree</code>	249	85	V_{pg}
<code>vm_kmap</code>	521	164	V_{kmap}
<code>vm_kunmap</code>	809	241	V_{kmap}
<code>vm_kmalloc</code>	908	285	V_{kmal}
<code>vm_kfree</code>	655	216	V_{kmal}
<code>vm_map</code>	-	-	V_{map}
<code>vm_unmap</code>	-	-	V_{map}

Hardware abstraction layer V_{hal} implements functions managing the virtual memory mappings and operating on MMU data structures.

Page allocator manages available physical memory using a page map structure. The map describes the accessible memory regions dividing it into segments consisting of page groups. Possible segment size values correspond to succeeding powers of 2 starting from 0. Each entry of the map points to a list of segments with a given size. The first entry points to the list of segments consisting of one page only. The N th entry points to the list of segments of 2^{N-1} pages. After the system initialization the memory is divided into segments with maximum sizes. During the allocation process, allocating function (`vm_pageAlloc`) looks for the first segment with a size greater or equal than requested one. When the segment size is greater than the requested size it is divided into two smaller segments. One of them is returned to the map. This process repeats until the obtained segment size value is closest the requested value. When segment is released (using function `vm_pageFree`) the kernel merges it with a succeeded or proceeded segment on the list. Described memory allocator (named buddy allocator) is typical for the most operating systems and has been presented in [9].

The kernel virtual space allocator V_{kmap} manages kernel virtual space V_K . An allocated segment is not accessible for the kernel and it should be mapped into the kernel address space if the kernel intends to use it. Mapping function `vm_kmap` looks for the first virtual space region at which a segment can be mapped. After mapping, the kernel marks this region as used. The last stage of mapping is the invocation of a hardware dependent `pmap_enter` function which fills structures used by MMU. Opposite function `vm_kunmap` removes segment from kernel virtual space, marks the region as free and inaccessible using `pmap_enter` function. Kernel space allocation is not a common mechanism in operating systems because many of them assume that the virtual address space is much greater than the physical memory size and the kernel space can correspond to the physical memory.

The kernel heap allocator is used to manage the allocation of memory chunks which size is less than the size of the page (minimal size of the segment). The main data structure of this allocator is the `size[]` array. Array entries point to segments divided into smaller chunks with specified size. The relation between the entry index (i) and the chunk size is given by the function 2^{4+i} . When the table entry is empty (no chunks are available) the new segment is allocated, mapped into the kernel address spaces and divided into chunks. Heap allocator consists of allocating function `vm_kmalloc` and releasing function `vm_kfree`. Each of these function depends on the segment allocation and mapping functions.

3 Error model and measures

The error model assumed in this paper is an error having impact on correct execution of an instruction by processor. This model is similar to the model presented in [5]. Single-bit errors are injected to have impact on instructions of functions belonging to the memory manager. Each injected error can be characterized by following attributes:

- *trigger* - An error is injected when instruction is fetched by processor.
- *location* - Errors are injected to code of selected functions at random locations. Errors are injected to any byte of an instruction.
- *type* - One single-bit (bit-flip) error per byte of an instruction is injected.
- *duration* - Two types of errors are considered - transient errors which are removed after the erroneous instruction execution and persistent errors which exist during the experiment execution time. Type of error is selected randomly, but the probability of selecting persistent error is 25%.

Outcomes from fault injection experiments are classified according to observable kernel states presented in table 2.

Table 2. Operating system states observable after the fault injection

State	Description
S_{err}	Injected fault has been detected using software mechanisms embedded into the kernel.
S_{exc}	Injected fault has been detected using CPU exception.
S_{hang}	Kernel hang ups and goes to a non-operational state. Restart is required.
S_{fs}	The corrupted instruction has been executed and no visible impact on the system has been observed either by the workload application and operating system kernel.
S_{fsv}	The corrupted instruction has been executed and requests generated by workload application are not properly serviced (i.e. two succeeded memory allocations returns the same heap chunk). This is the fail silence violation state.
S_{tm}	Injected error is not executed in the time assumed for the experiment.

To characterize what impact errors injected into the function f have on the operating system the special measure, called operating system resistance R , was introduced. Resistance R for function f is the vector of percentage shares of every observable state during the injection campaign for this function f (1).

$$R(f) = [p_{S_{err}}, p_{S_{exc}}, p_{S_{hang}}, p_{S_{fs}}, p_{S_{fsv}}, p_{S_{tm}}] \quad (1)$$

Other metrics characterizing error detection timing were introduced.

- average error activation time T_{act} —the time elapsed from the beginning of the experiment to the moment of execution of the corrupted instruction;
- standard deviation of the average error activation time σT_{act} ;
- average error detection time T_{det} —the time elapsed from the beginning of the experiment to the moment of error detection;
- standard deviation of average error detection time σT_{det} .

Error propagation is characterized for each function using a propagation ratio measure $Prop(f)$ (2).

$$Prop(f) = \frac{n_{prop}}{n_{act}} \quad (2)$$

Propagation ratio $Prop$ for the function f is the fraction of the number of errors detected outside the function (n_{prop}) to the number of all errors activated in this function (n_{act}). Value 0 of propagation ratio for function f indicates that all faults are detected during the function execution. Value 1 indicates that all faults are detected after passing control to other functions or to non-code segments.

Experiments has been performed for workload generated by thread calling routines used for validation of implementation of a memory management subsystem. These routines activate all functions of the memory manager.

4 Experimental setup

Dependability evaluation of the operating system requires different methodology than the dependability evaluation of applications running under the control of an operating system. Application running in presence of faults injected into its address space has no impact on operating system stability and all faulty conditions can be properly handled using debugging subsystem routines. Methodologies and tools developed for the dependability evaluation of user applications have been presented in [4].

Dependability evaluation of an operating system requires the existence of separate computer which controls fault execution and gathers outcomes of each experiment. This computer is called an experiment controller. Errors are injected by the fault injector embedded into an operating system kernel which communicates with the controller. Each fault injection corresponds to one experiment. Many experiments associated with a selected function constitute a campaign. The experiment process is presented below:

- **Select the location of the fault and other fault attributes**
- **Set the instruction breakpoint address**
Before setting the breakpoint instruction stream is disassembled to establish the location of the instruction containing the error location. Disassembled code is sent to the experiment controller. This step of the injection experiment is called the 'inject' phase.
- **Wait for the instruction execution**
When a selected instruction is fetched by the processor injector modifies the selected memory location and enables the step mode of instruction execution. This step is called the 'preinject' phase. Disassembled instruction stream and CPU context are sent to the controller.
- **Execute the corrupted instruction**
After the execution of the corrupted instruction the experiment enters into the 'postinject' phase. In this phase the injector sends CPU context and disassembled instruction stream starting from address of the program counter and restores the original value of the memory byte if the transient error is simulated.
- **Wait for the fault manifestation**
An experiment enters into the last stage which is the fault manifestation.
- **Restart the system using watchdog**

Outcomes from experiments are stored as files consisting data coming from each experiment state. The example file is presented below.

```
<inject time=12000 addr=c0060c90 bit=1 fault=t>
<code>
c0060c8e: movl 10(%ebx),%ebp
c0060c91: movl %ecx,ffffffc(%esi)
c0060c94: movl $fe8,%eax
</code>
</inject>

<preinject time=125000>
<cpu>
```

```

eax=ffe2de3b ebx=c018e000 ecx=c018e0cc
edx=c018e0a8 esi=c018e0ac edi=c006f12c ebp=0000b13c esp=00000001
eip=c0060c8e cs=00000008
ds=00000008 es=00000010 fs=00000010 gs=00000010 ss=00000010
err=00000000 eflags=00000186
</cpu>

<code>
c0060c8e: movl 12(%ebx),%ebp
c0060c91: movl %ecx,ffffffc(%esi)
c0060c94: movl $fe8,%eax
c0060c99: movl c(%ebx),%esi
</code>
</preinject>

<postinject time=125000>
<cpu> eax=ffe2de3b ebx=c018e000
ecx=c018e0cc edx=c018e0a8 esi=c018e0ac edi=c006f12c ebp=14b40000
esp=00000001 eip=c0060c91 cs=00000008
ds=00000008 es=00000010 fs=00000010 gs=00000010 ss=00000010
err=00000000 eflags=00000086
</cpu>

<code>
c0060c91: movl %ecx,ffffffc(%esi)
c0060c94: movl $fe8,%eax
c0060c99: movl c(%ebx),%esi
c0060c9c: movl %edx,8(%ebx)
</code>
</postinject>

<timeout time=60460800> </timeout>
    
```

This file describes the fault injected 12 ms after the system initialization. This injection modifies the second bit of byte belonging to the `movl` instruction at address `0xc0060c8e`. The corrupted instruction has been executed 125 ms after the system initialization. Fault modifies offset of the operand moved to the EBP register. After the instruction execution the EBP register stored incorrect value but this error had no impact on workload and the operating system (fail silence state S_{fs}). After the timeout system was restarted and the next experiment was performed.

5 Experimental results

During the memory subsystem tests, 1503 experiments have been performed. Outcomes from experiments are divided into two groups corresponding persistent and transient errors. For each function of the memory subsystem introduced measures are calculated. Table 3 presents operating system resistances R for persistent errors injected into the memory management functions. Table 4 presents resistances for transient errors.

Table 3. Operating system resistance for persistent error injected into the memory manager

Function	N	ps_{err}	ps_{exc}	ps_{hang}	ps_{fs}	ps_{fsv}	ps_{tm}
<code>pmmap_enter</code>	124	0% (0)	1% (1)	46% (57)	10% (13)	0% (0)	43% (53)
<code>vm_pageAlloc</code>	46	0% (0)	0% (0)	13% (6)	30% (14)	0% (0)	57% (26)
<code>vm_pageFree</code>	41	0% (0)	5% (2)	27% (11)	58% (24)	0% (0)	10% (4)
<code>vm_kmap</code>	27	0% (0)	4% (1)	55% (15)	4% (1)	0% (0)	37% (10)
<code>vm_kunmap</code>	36	0% (0)	3% (1)	33% (12)	6% (2)	0% (0)	58% (21)
<code>vm_kmalloc</code>	97	0% (0)	13% (13)	3% (3)	21% (20)	0% (0)	63% (61)
<code>vm_kfree</code>	39	0% (0)	10% (4)	8% (3)	38% (15)	0% (0)	44% (17)

The column named N presents the number of experiments conducted during the campaign. The number of experiments is relatively low in comparison to the number of function instructions but the

goal of experiments was not to evaluate dependability of the memory manager in details . The goal of the experiments was to present the methodology and framework for a quick dependability estimation of implemented kernel functions using introduced measures. Such estimation allows for selecting functions requiring introduction of propagation barriers or enhancing error detection mechanisms. From development process perspective it makes no sense to evaluate in details dependability of every function because implementation is changing during the system life-cycle. Other columns present items of the resistance vector.

Table 4. Operating system resistance for transient error injected into the memory manager

Function	N	ps_{err}	ps_{exc}	ps_{hang}	ps_{fs}	ps_{fsv}	ps_{tm}
<code>pmap_enter</code>	367	0% (0)	4% (16)	50% (181)	9% (33)	0% (0)	37% (137)
<code>vm_pageAlloc</code>	145	0% (0)	2% (3)	14% (20)	25% (36)	0% (0)	59% (86)
<code>vm_pageFree</code>	101	0% (0)	3% (4)	31% (31)	52% (53)	0% (0)	13% (13)
<code>vm_kmap</code>	87	0% (0)	2% (2)	40% (35)	11% (9)	0% (0)	47% (41)
<code>vm_kunmap</code>	120	0% (0)	3% (3)	22% (26)	5% (6)	0% (0)	70% (85)
<code>vm_kmalloc</code>	108	0% (0)	13% (14)	1% (1)	22% (24)	0% (0)	63% (69)
<code>vm_kfree</code>	162	0% (0)	7% (11)	7% (12)	28% (45)	0% (0)	58% (94)

Results show that many of injected faults are not executed. Most faults are ignored in functions `vm_kmalloc` and `vm_kunmap`. This is a consequence of workload nature. Typical execution of `vm_kmalloc` requires no page allocation and the branch instruction at the function beginning omits execution of a large part of code. Similar scenario takes place in `vm_kunmap`. When a kernel address space region is released, the kernel tries to concatenate it with other free regions on the list. If no candidate for concatenation is found, branch instruction moves execution to the end of the function. Least faults are omitted in function `vm_pageFree`. The function code contains many branches but they omit only short parts of the code.

The other important fact should be noted. Results show that software error detection mechanisms like assertion verifying correctness of input arguments passed to a function are insufficient. No error has been detected using this mechanism. Paper [6] is focused only on injecting faults into input arguments of functions and his author concluded that wrapping is a good technique for prevention of fail silence violations and increasing the number of detected errors. Authors assumed that measures introduced by them could be helpful in finding vulnerable operating system parts which potentially demand enhancing argument validation. Presented results show that such methodology is insufficient. The percentage of errors detected using assertion and exceptions is extremely low when faults are injected into the base system component which is crucial for overall system stability.

In presented results there are no outcomes associated with fail silence violation state S_{fsv} . This is very difficult to verify outputs of memory allocation and mapping functions. If the mapping function fails and validating function tries to reference the mapped segment system hangs or processor raises the exception. This is the possible explanation of such situation. This explanation is convinced by the relatively high number of system hangups. Most hangups were observed when faults were injected into mapping function which operate directly on MMU.

Table 5. Detection times and propagation ratio for persistent errors injected into the memory manager

Function	T_{act} [μs]	σT_{act} [μs]	T_{det} [μs]	σT_{det} [μs]	Prop
<code>pmap_enter</code>	114027	3616	111000	0	0.00000
<code>vm_pageAlloc</code>	137379	6869	-	-	0.00000
<code>vm_pageFree</code>	146455	5595	1862500	2413355	0.00000
<code>vm_kmap</code>	109600	1131	108800	0	0.00000
<code>vm_kunmap</code>	114000	6245	109000	0	0.00000
<code>vm_kmalloc</code>	116829	5537	116692	5212	0.05556
<code>vm_kfree</code>	115200	3211	115200	4460	0.00000

Tables 5, 6 present detection times and propagation ratio for each tested function. Injected errors have been detected practically after activation. In some cases only the latency between activation and detection

was about 1 s. Highest error propagation was observed for `vm_kmalloc` and `vm_kunmap` functions. This is the consequence of function complexities and dependencies with other memory manager functions.

Presented results for persistent and transient errors are similar. When fault is injected it is practically manifested after the first execution of corrupted instruction. This conclusion is convinced by the measured short times between error activation end error detection.

Table 6. Detection times and propagation ratio for transient errors injected into the memory manager

Function	T_{act} [μs]	σT_{act} [μs]	T_{det} [μs]	σT_{det} [μs]	$Prop$
<code>pmap_enter</code>	113100	3799	111787	2901	0.00435
<code>vm_pageAlloc</code>	135396	5429	129933	2003	0.00000
<code>vm_pageFree</code>	145813	4919	538950	787508	0.02273
<code>vm_kmap</code>	111691	4461	110200	3959	0.02174
<code>vm_kunmap</code>	113556	4086	113800	4275	0.05714
<code>vm_kmalloc</code>	114958	4261	814100	2611190	0.05128
<code>vm_kfree</code>	116035	3739	115615	4768	0.01471

6 Conclusions

The paper presents methodology for dependability evaluation of Phoenix-RTOS real-time operating system developed by author for real application needs. Dependability evaluation is the last stage of implementation of every Phoenix-RTOS function. Outcomes from fault injection experiments allow to determine the kernel resistance (efficiency of error detection mechanisms) for errors introduced into the selected function and the function ability to propagate errors. In opposite to popular real-time operating system, which authors concentrate only on system stability, this is the novel approach which allows to propose new kernel-level error detection mechanisms and propagations barriers. Obtained results related to the memory manager (presented in section 5) allowed to select functions demanding the enhancement of error detection mechanism and introduction of propagation barriers. According to proposed methodology the goal of the introduction of error detection enhancements and error propagation barriers is to obtain (for every VM function) specific error resistance vectors with dominant $p_{S_{exc}}$ and $p_{S_{err}}$ components.

Furthermore obtained results show that techniques based on wrapper, proposed in the paper [6]), are insufficient when errors are injected into the functions of crucial operating system parts.

Obtained results will be used in next stage of Phoenix-RTOS development to develop the dependable memory manager. The last stage of the development process will be the dependability evaluation of the health monitoring system based on Phoenix-RTOS.

Acknowledgment

This work was supported by KBN grant 4T11C049 25.

References

- Ademaj A.: *A Methodology for Dependability Evaluation of the Time-Triggered Architecture Using Software Implemented Fault Injection*, Fourth European Dependable Computing Conference (EDCC-4) 2002.
- Arlat J., Fabre J. C., Rodriguez M., Salles F.: *Dependability of COTS Microkernel-Based Systems*, IEEE Transaction on Computers **Vol. 51**, No. 2 2002.
- Fuchs E.: *An Evaluation of the Error Detection Mechanisms in MARS using Software-Implemented Fault Injection*, 2nd European Dependable Computing Conference (EDCC-2) 1996.
- Gawkowski P.: *Analysing an enhancing fault immunity of programs in systems with COTS elements*, PhD. thesis Institute of Computer Science Warsaw University of Technology 2005.
- Weining Gu, Kalbarczyk Z., Iyer Ravishankar K., and Zhenyu Yang: *Characterization of Linux Kernel Behavior under Errors*, IEEE DSN 2003.
- Johansson A., Neeraj Suri: *Error Propagation Profiling of Operating Systems*, IEEE DSN 2005.
- Pisarczyk P.: *Phoenix—realtime kernel for embedded applications*, Real-Time Systems Conference Ustron, 2002.
- Rodriguez M., Salles F., Fabre J. C., Arlat J.: *MAFALDA: Microkernel Assessment by Fault Injection and Design Aid*, 3rd European Dependable Computing Conference (EDCC-3) 1999.
- Tanenbaum, A. S.: *Modern Operating Systems*, 2nd Edition Prentice-Hall, 2001.