

# In-System Programming of Non-Volatile Memories on Microprocessor-centric Boards

Anton Tsertov, Sergei Devadze, Artur Jutman, and Artjom Jasnetski

**Abstract**—With the continuous growth of capacity of non-volatile memories (NVM) in-system programming (ISP) has become the most time-consuming step in post-assembly phase of board manufacturing. This paper presents a method to assess ISP solutions for on-chip and on-board NVMs. The major contribution of the approach is the formal basis for evaluation of the state-of-the-art ISP solutions. The proposed comparison pin-points the time losses, that can be eliminated by the use of multiple page buffers. The technique has proven to achieve exceptionally short programming time, which is close to the operational speed limit of modern NVMs. The method is based on the ubiquitous JTAG access bus which makes it applicable for the most board manufacturing strategies despite a slow nature of JTAG bus.

**Index Terms**—in-system programming, processor-centric board, JTAG, non-volatile memory.

## I. INTRODUCTION

THE widely adopted DfT structures defined in IEEE 1149.1 standard [1] and complemented in [2] are heavily used for ISP and memory test, besides traditional post-assembly tests (Boundary-Scan tests). Despite of ubiquitous presence of boundary-scan (BS)[1] structures in modern electronic systems and components, the application of BS is limited due to the low operation frequency. Typically BS clock (TCK) frequency is in range from 1 MHz to 50 MHz, whereas actual speed of data transfer is much lower because of the overhead data (JTAG instructions and UUT protocol) that accompanies each test pattern. Here and later the term BS will be used to stress that IEEE 1149.1 structures are used in isolation from the internal functionality of the chip. The term JTAG denotes the test infrastructures defined in IEEE 1149.1 standard.

Traditionally, the ISP is considered to be the final phase of BS test session and it is essential for the functional tests [3]. The fundamental JTAG-based ISP solutions are proposed in [4][5] and [6], that describe processor-controlled test and processor-centric board test (PCBT). These methods are common in the use of functionality of microprocessor (uP) or microcontroller (uC) to access peripheral components outside the uP or uC at normal operational speed of the board. The attractiveness of such uP-based board test solutions are very high due to its non-intrusiveness into board or/and chip design.

A. Tsertov is with Tallinn University of Technology, Department of Computer Engineering, Akadeemia 15a, 12618 Tallinn, Estonia (anton.tsertov@ttu.ee).

S. Devadze, A. Jutman and A. Jasnetski are with Testonica Lab, Akadeemia 15a, 12618 Tallinn, Estonia (sergey@testonica.com, artur@testonica.com, artjom@testonica.com).

Later in this paper, the formal characterization of ISP will be given in the terms from JTAG and PCBT methodology.

The ISP of NVMs with the use of PCBT methodology has proven to speed up the existing JTAG-based solutions [5], [7]. However, the size of the system firmware (program image) that is stored to the NVMs is following the trend that is set by continuously growing hardware functionality. Inevitably the speed up in ISP time introduced by PCBT methodology is not sufficient and the industry demands ISP solutions that will run at the operational speed limit of contemporary and future NVMs. State-of-the-art JTAG and UART-based solutions are limited by the speed of the data transfer link. In this paper we derive the technique that mitigates time-losses that are common to most ISP approaches.

### A. Problem statement and paper contribution

After the bare PCB is populated with components it needs to be tested for manufacturing defects. Typically, the boundary scan tests are used to screen out boards with static structural defects. After the screening procedure the “healthy” boards reside in the fixture for subsequent ISP and functional testing, e.g. with the use boot-loader image. In most cases, it is considered beneficial to program NVMs with the same tester hardware that is used for verification and test of other components on the printed circuit board assembly (PCBA) under test. The PCBT methodology helps to reduce the programming time of flash memory from hours, as in case with BS, to minutes and even less. Nevertheless, the actual ISP time (in case of PCBT) heavily depends on the architecture of the debug interface of the uP, on the instruction set of the uP, on the performance of the flash memory controller inside the uP SoC and on the performance of the flash memory itself. The latter is discussed in details in the last section of this paper.

The PCBT-based ISP procedure of NVMs is formally characterized in Section II. Section III studies the drawbacks of the state-of-the-art ISP solution and proposes countermeasures to solve the named issues. Section IV is devoted to the approach that allows in most cases to perform in-system programming of on-chip or on-board non-volatile memories at maximum speed. Whereas, maximum speed means that the bottleneck is not in the data transfer channel (JTAG), but in the capability of the memory itself to program the supplied data faster. The formal basis from Section II and III is used in the last section to assess and compare the experimental results of different test-cases (ISP of on-chip and on-board NVM).

## II. IN-SYSTEM PROGRAMMING

The in-system programming is a process of loading data or program image into memory. The memory can be either the on-board non-volatile memory or the on-chip (e.g. embedded NVM of the uP). The *in-system* key-word means that the memory resides inside the system while it is being programmed.

In general, the ISP of NVMs consists of the following steps:

1. *Unlock* - Unlock NVM for write operations
2. *Erase* - Do erase for the area to be programmed
3. *Blank check* - Check the erased area to be blank
4. *Program* - Program data to the NVM
5. *Verify* - Verify the programmed data

Traditional BS-based approach is capable to run all these steps, in case the connections between NVM and the BS-enabled device are controllable from the BS chain. The length of the BS chain is often measured in thousands of bits in modern systems, hence the useful data payload in BS-based communication is extremely small, especially for memories with serial interface. The next section shows how ISP can benefit from the use of uP functionality instead of BS register to speed up programing of NVMs.

### A. Processor-centric Board Test

The processor-centric board test [6] is a collective term for the post-manufacturing tests for processor-centric boards. These tests target the defects related to the last stages of the board manufacturing process and also the first period of the post-manufacturing product life (e.g. in-system programming (ISP), infant mortality diagnosis). The uP plays a role of on-board tester, that listens to commands from external test equipment and in response applies tests to other PCBA components. The details on implementing ISP with the use of PCBT methodology are described in the following paragraphs.

Test path initialization and configuration belong to the **test access** functionality of the PCBT program. The rest of the PCBT program functionality is a part of the **test application**, which may be developed in accordance to online<sup>1</sup> or offline test application modes [8].

In case of the offline (autonomous) mode, the complete test program (test vectors and expected values) is translated into the set of uP instructions and loaded as an ordinary program into embedded memory of the uP. The program execution on the uP is started by the external tester. After test program execution is finished the result (PASS or FAIL for complete test) will be stored in the on-chip memory of the uP. It is retrieved through the debug interface and reported to the external tester for further evaluation and diagnosis.

The offline mode is fully independent and does not suppose continuous interaction with an external tester. This mode needs available on-chip memory to store test data. Obviously, this mode is not suitable for ISP of large images unless there is enough memory available (e.g. on-chip or on-board volatile memory) to hold the whole image to be programmed into on-board non-volatile memory.

The key difference between the online and offline modes is that in the online mode commands are executed under the control of external tester. In online mode the flash image is transferred through the test access path word by word directly to the non-volatile on-board memory. This implies also transferring the commands for NVM, which significantly reduces the test data throughput of test application.

### B. Calculation of ISP time

In this section the formal basis for ISP time estimation is proposed. Each ISP step will be characterized in PCBT terms with the respective formula.

1) *Programming*: The formula 1 is proposed for time calculation for programming step ( $t_{ISPp}$ ) in case of online mode of test application. It has summands for shifting in data image, the uP instructions for writing data image and for reading and shifting out the NVM status after buffer program operation.

$$t_{ISPp} = t_{WD} + dt_{IW} + t_B \frac{d}{b} + t_{IR} \frac{d}{b} + \frac{t_{RD}}{b} \quad (1)$$

The first two summands in formula 1 is for writing data to the pages and the last two summands are for checking the result of page program operation.

- $d$  - size in words of the flash image
- $b$  - size in words of the flash page. Traditionally the data in NVMs is handled page-wise.
- $t_{WD}$  - time to transfer over JTAG bus (shift in) the flash image  
 $t_{WD} = d \frac{l}{h}$ , where  $l$  is a length of the shift and  $h$  is a frequency of the test clock (TCK).
- $t_{RD}$  - time to transfer over JTAG bus (shift out) the flash image  
 $t_{RD} = d \frac{l}{h}$ , where  $l$  is a length of the shift and  $h$  is a frequency of the test clock (TCK). One can notice that  $t_{RD}$  and  $t_{WD}$  are identical, hence, for simplification, in the following discussions  $t_{WD}$  is used instead of  $t_{RD}$ .
- $t_{IW}$  - time to shift in the instructions that write a word to any memory location  
 $t_{IW} = \sum_{i=1}^k \frac{l_i}{h}$ , where  $k$  is the number of shifts to emulate the write instruction from external tester on the uP.  $k$  includes the JTAG TAP instructions, debug port instructions and uP instructions.  $l_i$  denotes that instruction of every type might be different in length.
- $t_B$  - programming time, required by Flash device to program one page  
 $t_B$  - is a constant specified in the datasheet of the memory
- $t_{IR}$  - same as  $t_{IW}$ , but for read  
 $t_{IR} = \sum_{i=1}^r \frac{l_i}{h}$  where  $r$  is the number of shifts to emulate the read instruction from external tester on the uP.  $r$  includes the JTAG TAP instructions, debug port instructions and uP instructions.  $l_i$  denotes that instruction of every type might be different in length.

$\frac{d}{b}$  means the number of pages to program (in case of unaligned number of bytes the result is rounded up).

<sup>1</sup>Online mode of test application is not the same as on-line testing

2) *Verify and Blank Check*: The verification of the programmed data can be done in two ways. First way is to read data from the NVM to external tester and verify it against the programmed data externally. This is suitable for online mode of PCBT. Second way is to read data from NVM to on-board volatile memory (e.g. uP embedded RAM or on-board RAM/DRAM) and let the uP to execute the comparison routine. The last one is beneficial for offline mode.

The time for verification ( $t_{ISP_V}$ ) of programmed data in online mode is proposed to estimate with formula 2.

$$t_{ISP_V} = t_{RD} + dt_{IR} + t_{IR} \frac{d}{b} + \frac{t_{RD}}{b} \quad (2)$$

The notations in formula 2 have the same meaning as in formula 1. The first two summands stand for time that is required to read the programmed data back to external tester. While the last two summands show the time that is needed to read the status of the NVM in order to check for errors that may appear during the page read operation.

For *Blank Check* operation the content of NVM can be checked either inside the uP or in external tester. The check inside the uP requires additional software routine which is not common for offline mode. In online mode the verification routine can be used to verify the read data against blank value (e.g. *FFFFh*).

3) *Unlock and Erase*: The *Unlock* and *Erase* operations consist only from respective NVM command and an additional commands to check the status register. In formula 3 (*Unlock*) and 4 (*Erase*) commands are used for each page, this adds redundant operations for memories that allow block-wise (multi-page) unlock and erase operations. However, this allows same formulas to be applicable for all NVM types and from the functional perspective given redundancy does not cause malfunctions.

$$t_{ISP_U} = (t_{IWu} + t_{IR} + t_{UN}) \frac{d}{b} \quad (3)$$

- $u$  - number of commands to unlock a page
- $t_{UN}$  - time required by NVM to unlock one page

$$t_{ISP_E} = (t_{IWc} + t_{IR} + t_{ER}) \frac{d}{b} \quad (4)$$

- $c$  - number of commands to erase a page
- $t_{ER}$  - time required by NVM to erase one page

### III. HYBRID ONLINE ISP MODE

In order to speed up the ISP in online mode the industry has came up with a solution [5], [9], which we call in the following as hybrid-online mode of test application. In Figure 1 the uP-based SoC components that are active in hybrid mode are shown. The difference to pure online mode is in the software that is executed by the uP. This software listens to the control commands from external tester and is called *monitor software*.

#### A. Programming

Firstly, the external tester transfers part of the data image to the buffer, then the command to copy data from buffer to buffer inside NVM is send to the monitor software. The monitor software handles the data transfer from the buffer to the page buffer inside NVM by turn and sends write page commands to NVM. If the on-chip memory is not available the on-board volatile memory can be used to store the monitor software and host the intermediate buffer for data. The NVM programming time in hybrid-online mode is proposed to calculate with formula 5.

$$t_{ISP_P} = t_{WM} + t_{WD} + (m + d)t_{IW} + t_B \frac{d}{b} + t_{EX_P} \frac{d}{b} \quad (5)$$

- $m$  - size in words of the monitor software
- $t_{WM}$  - time to shift in the monitor software  
 $t_{WM} = m \frac{l}{h}$ , where  $m$  is size of the monitor in words,  $l$  is a length of a word and  $h$  is a frequency of the test clock (TCK).
- $t_{WD}$  - time to shift in the flash image  
 $t_{WD} = d \frac{l}{h}$ , where  $l$  is a length of a word and  $h$  is a frequency of the test clock (TCK).
- $t_{IW}$  - time to shift in instructions that write a word to the register or location in volatile memory  
 $t_{IW} = \sum_{i=1}^k \frac{l_i}{h}$ , where  $k$  is the number of shifts to emulate the write instruction from external tester on the uP.  $k$  includes the JTAG TAP instructions, debug port instructions and uP instructions.  $l_i$  denotes that instruction of every type might be different in length.
- $t_B$  - programming time per page  
 $t_B$  - is a constant specified in the datasheet of the memory
- $t_{EX}$  - time taken by monitor to copy flash image from buffer to page buffer into flash memory. As the monitor is executed at the actual operating speed of the uP this time can be neglected (proven in Section V). Typically uP clock is at least one order of magnitude faster than TCK.

After substituting the notation of summands in formula II with the respective expressions, formula 5 takes the following form:

$$t_{ISP} = m \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) + d \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) + t_B \frac{d}{b} \quad (6)$$

Obviously, in order to justify the effort spent on development of the monitor software the hybrid-online mode has to speed up the online mode. This consideration is expressed by the following inequality  $1 > 5$ :

$$\frac{d}{h} + dt_{IW} + t_B \frac{d}{b} + \frac{d}{b} t_{IR} > m \frac{l}{h} + d \frac{l}{h} + (m + d)t_{IW} + t_B \frac{d}{b}$$

that reduces to:

$$\frac{d}{b} t_{IR} > m \frac{l}{h} + mt_{IW} \quad (7)$$

Let us evaluate the obtained expression 7. From 7 it can be concluded that the time to load the monitor to volatile

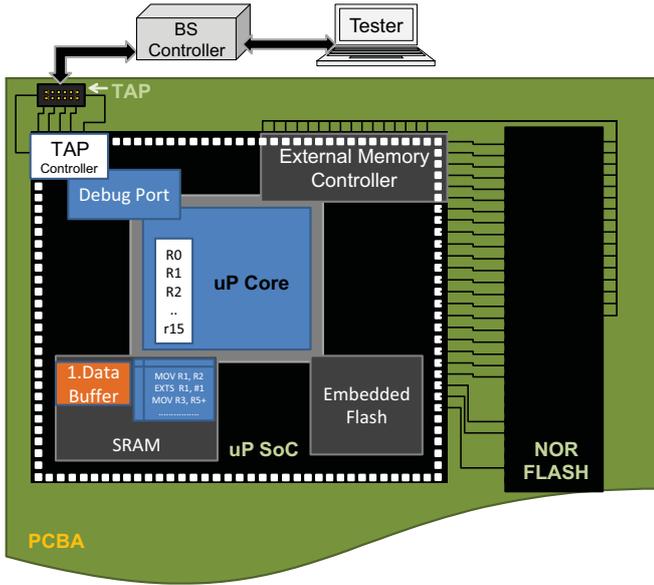


Fig. 1. Data path components in hybrid online mode of ISP

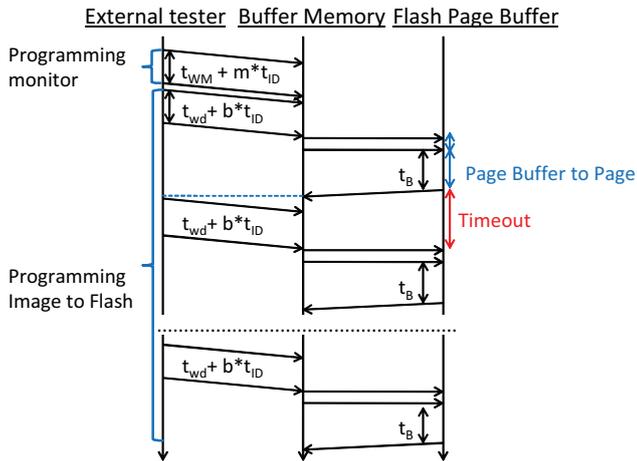


Fig. 2. Time flow in hybrid online mode of ISP

memory should be shorter than the time taken by reading the flash status after the page programming operations. In other words, the inequality 7 never holds if the monitor size in words is bigger than the number of polls for flash status (assuming that the read operation and write operation takes the same number of shifts). Number of polls for flash status depends on the number of flash pages to program, which has direct relation to the size of the image to be programmed. Hence, the programming operation in hybrid-online mode (in the form it is described here) will likely to be slower than the online mode for small images.

### B. Verify

Standalone verification of programmed data in online mode is performed inside the uP. Hence it requires transfer of original data to the volatile memory. The time for data transfer is reflected by the first two summands in equation 8. The  $t_{EX_v}$  stands for time that monitor uses to read the data back from

flash pages and compare it against the original data. Normally, the verification is done page-wise and the flow-graph is similar to programming in Figure ??monitor-ISP. It is assumed that the monitor program is loaded to uP program memory in the previous steps.

$$t_{ISP_V} = t_{WD} + dt_{IW} + t_{EX_V} \frac{d}{b} \quad (8)$$

The standalone verification is almost as much time costly as a programming step itself. The verification of programmed data can be joined with programming step. After the page data is programmed to the NVM page the original image data for that page still resides in the volatile memory. Thus, page verification after page programming step saves time for transferring data from external tester. Then the time for verification consists only from  $t_{ISP_V} = t_{EX_V}$ .

### C. Unlock, Erase, Blank Check

In hybrid online mode the *Unlock*, *Erase* and *Blank Check* are implemented inside the *monitor software*. The external tester tells the monitor the start page, the number of pages and the step to execute (unlock, erase or blank check). Altogether it forms three command words to send to the monitor -  $3(t_{IW} + \frac{l}{h})$ . The formulas 9, 10, 11 are proposed to calculate time taken by unlock, erase and blank check steps of ISP in hybrid online mode.

$$t_{ISP_U} = 3(t_{IW} + \frac{l}{h}) + \frac{d}{b}(t_{EX_{UN}} + t_{UN}) \quad (9)$$

$$t_{ISP_E} = 3(t_{IW} + \frac{l}{h}) + \frac{d}{b}(t_{EX_{ER}} + t_{ER}) \quad (10)$$

$$t_{ISP_{BC}} = 3(t_{IW} + \frac{l}{h}) + \frac{d}{b}t_{EX_{BC}} \quad (11)$$

- $t_{EX_{UN}}, t_{EX_{ER}}, t_{EX_{BC}}$  - time that monitor software uses to execute the respective command
- $t_{UN}$  - time that NVM needs to unlock one page
- $t_{ER}$  - time that NVM needs to erase one page

### D. Summary

Let us join all ISP steps in one equation:

$$\begin{aligned} t_{ISP} &= t_{WM} + t_{WD} + (m + d)t_{IW} + t_B \frac{d}{b} \\ &\quad + t_{EX_P} \frac{d}{b} + t_{EX_V} \frac{d}{b} + 9(t_{IW} + \frac{l}{h}) \\ &\quad + \frac{d}{b}(t_{EX_{ER}} + t_{ER} + t_{EX_{UN}} + t_{UN} + t_{EX_{BC}}) \\ &= t_{WM} + t_{WD} + (m + d)t_{IW} \\ &\quad + 9(t_{IW} + \frac{l}{h}) + \frac{d}{b}(t_{UN} + t_{ER} + t_B) \\ &\quad + \frac{d}{b}(t_{EX_{UN}} + t_{EX_{ER}} + t_{EX_{BC}} + t_{EX_P} + t_{EX_V}) \end{aligned} \quad (12)$$

Assuming that normally number of words to program to NVM is exceeding hundreds of thousands we may neglect the time that is used to transfer commands for unlock, erase

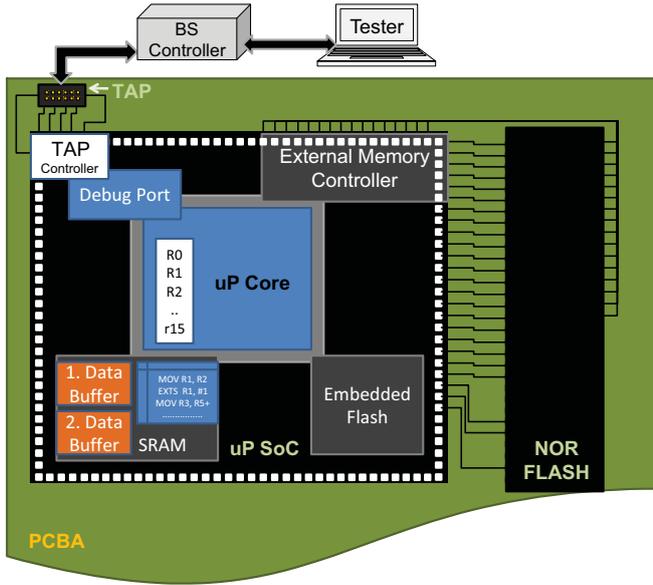


Fig. 3. Data path components in double-buffer hybrid online mode of ISP

and blank check operations, as well as to neglect the monitor execution time for all operations. The monitor execution time is calculated in microseconds even for slow uP, while our calculations in the following sections are in milliseconds. Then, considering the simplification proposed here the formula 12 takes the following form:

$$t_{ISP} = t_{WM} + t_{WD} + (m + d)t_{IW} + \frac{d}{b}(t_{UN} + t_{ER} + t_B) \quad (13)$$

#### IV. DOUBLE BUFFER ONLINE ISP MODE

Despite the hybrid-online mode in general gives shorter ISP time than in online mode, there is still a place for optimizations. It should be stressed that in the ideal ISP solution the bottleneck is in the time taken by the NVM itself to program the page ( $t_B$ ). In figure 2 is shown the time flow in hybrid-online mode of ISP. In other words, in ideal solution ISP is as fast as the flash can allow. Hence, the timeout shown in Figure 2 needs to be minimized. The time that is used to program data from volatile buffer to page buffer inside flash ("Buffer to Page Buffer") and time  $t_B$  (Page Buffer to Page) are inevitable in any case.

In order to catch up with the time of the ideal solution let us consider adding another buffer to the hybrid-online mode of ISP. The idea is to use the time that flash needs to program the data in its page buffer to the actual page ( $t_B$ ) for transferring data for the next page from external tester to the buffer in volatile memory. The data flow for double-buffered hybrid online method is shown in Figure 4. In this case the programming of the flash page and the data transfer via test access path are performed in parallel. The limitation of this approach is that the volatile memory has to be big enough to store the monitor and two pages (Buffer1 and Buffer2) of flash image (see Figure 3).

The formula 14 is proposed to calculate the ISP time for the double-buffer approach. The formula 14 is derived from 6 with

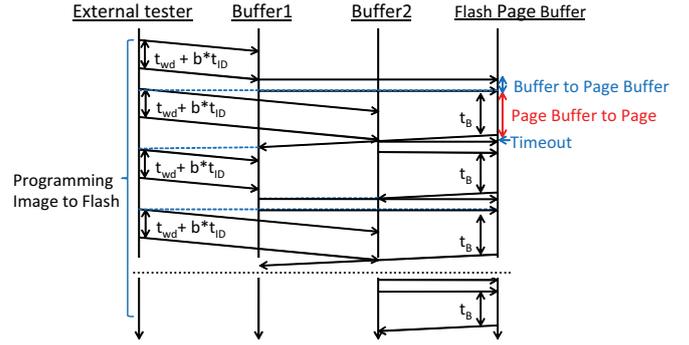


Fig. 4. Time flow in double-buffer hybrid online mode of ISP

assumption 15. In 15 is assumed that buffer programming time is shorter than a flash page programming time ( $t_B$ ), which is given in the datasheet. As it is shown later the assumption 15 holds in case of most on-chip and on-board flashes unless the flash has extraordinary small page programming time. In addition, formula 14 bounds the speed of communication with external tester ( $h$ ) from the lower side, when the equation is satisfied.

$$t_{ISP_P} = m \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) + b \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) + t_B \frac{d}{b} \quad (14)$$

$$b \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) \leq t_B \quad (15)$$

In formula 15 is shown that ISP time for programming ( $t_{ISP_P}$ ) is formed by the monitor programming to volatile memory, transferring the first page to buffer inside volatile memory (the rest of the pages are transferred in parallel to the programming of the previous page to flash) and the page programming time multiplied by the number of pages to program.

The applicability of the double-buffer hybrid-online mode of ISP in the particular test case can be evaluated using inequality 15. When the inequality is satisfied the overall ISP time is limited by the flash performance, otherwise the double-buffer based approach does not give a significant speed up and the ISP is limited by the throughput of the test access path.

#### V. FROM THEORY TO PRACTICE

The efficiency of proposed methodology is studied for two use-cases. The first use-case considers programming of embedded (on-chip) flash of the uC. In the second use-case the external (on-board) flash is programmed.

##### A. On-chip NVM

For the first use-case the uC XC2361E from Infineon [10] is used. The embedded flash memory of XC2361E is built from pages of size 128 bytes. The typical programming time for single page is 3ms [10]. Thus, the theoretical programming time of 256KB is 6144ms. The experimental results presented in Table I show the difference in programming 256KB of data into embedded non-volatile memory using toolchain from Keil [11], PCBT approach and modified PCBT approach that is

TABLE I  
PROGRAMMING OF 256KB INTO INFINEON XC2300E uP ON-CHIP FLASH

Method	Program Time (s)	Throughput in KB/s
Keil (UART)[11]	151,38	6,76
Keil (JTAG)[11]	14,72	17,12
PCBT [12]	12,43	20,27
PCBT (2-buff.)	7,81	32,25
Theoretical	6,144	35,714

based on the proposed double buffer technique. PCBT approaches are executed on the toolchain from Goepel Electronic [12]. The results from Keil-based toolchain show the state-of-the-art approach time. It should be mentioned that the actual communication frequency is not available for that approach, while for the PCBT the 20MHz test clock (TCK) frequency was used.

The double buffer PCBT approach outperformed other approaches and showed the time close to the theoretical one. The time for experiments listed in Table I include not only the ISP time, but also the board and uP initialization time. Hence, the difference between theoretical estimation and the double buffer PCBT approach represents the time for initializing/booting the board and uP and also the time for programming monitor software and the first page into intermediate buffer:  $7810 - 6144 = 1666(ms)$ . The pure programming time for hybrid online PCBT is  $12430 - 1666 = 10764(ms)$ . Hence, the time for transferring 256KB of data to intermediate buffer is  $10764 - 6144 = 4620(ms)$  and time to transfer 1 page to that buffer is  $4620/(2048 - 1) = 2.26(ms)$ . This shows that inequality 15 holds, whereas the left part is equal to  $2.26ms$ . and the  $t_B = 3ms$  as was stated previously.

### B. On-board NVM

For the second experiment the board with Emerald-P microprocessor and NOR flash from Numonix was selected. Emerald-P is based on ARM Cortex-A9 core that is paired with ADIV5 debug interface from ARM. The experimental results for ISP of 1 MB image into on-board NOR flash and the theoretical expectations both for hybrid online-hybrid PCBT ( $ISP_{OH}$ ) and two-buffer online-hybrid PCBT ( $ISP_{2BOH}$ ) approaches are shown in Table II. The theoretical result is based on the data from respective Numonyx manual [13] for selected NOR flash, where the programming time is defined as a range ( $min < t_B < max$ ) for selected 512 byte long page buffer. From Table II it is seen that the experimental result for double buffer PCBT does not fit into the theoretical range. Which may mislead us to the conclusion that for the given case the bottleneck is in the throughput of test access path, regardless the noticeable time improvement achieved in PCBT conditioned by implementation of the second buffer.

To discover the actual bottleneck let us find out the exact value of  $t_B$  parameter and the experimentally achieved throughput value.

In this experiment the size of the flash page and the buffer size is 512 bytes. From data shown in Table II the overall time to program one flash page in PCBT is  $t_{ISP_{P-OH}} =$

$(4434 - 100)/2048 = 2.12(ms)$ , where 2048 is the number of buffers in 1 MB. The same operation for the double buffer PCBT takes:  $t_{ISP_{P-2BOH}} = (2930 - 100)/2048 = 1.38(ms)$ . The difference between these two approaches conforms to  $t_{ISP_{P-OH}} - t_{ISP_{P-2BOH}} = 4434 - 2930 = 1504(ms)$ . The same in the formal view:

$$\begin{aligned} & t_{ISP_{P-OH}} - t_{ISP_{P-2BOH}} \\ &= d \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) + t_B \frac{d}{b} - b \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) - t_B \frac{d}{b} \\ &= (d - b) \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) \end{aligned} \quad (16)$$

If

$$b \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) \leq t_B$$

Otherwise:

$$\begin{aligned} & t_{ISP_{P-OH}} - t_{ISP_{P-2BOH}} \\ &= d \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) + t_B \frac{d}{b} - d \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) - t_B \\ &= t_B \left( \frac{d}{b} - 1 \right) \end{aligned} \quad (17)$$

In order to state whether inequality 15 holds or not it is needed to calculate the exact  $t_B$  value. Let us pessimistically assume that inequality 15 does not hold, then from 17  $t_B = 1504/(d/b - 1) = 1504/2047 = 0.73(ms)$ .

The experiments are executed using 20MHz clock for TCK signal ( $h = 20000(ms)$ ). Thus, the inequality 15 takes the following form after substituting parameters with their values:

$$512 \left( \frac{32}{20000} + \sum_{i=1}^k \frac{l_i}{h} \right) \leq 0.73$$

after simplifying:

$$0.82 + 512 \sum_{i=1}^k \frac{l_i}{h} \leq 0.73$$

Which indeed leads to contradiction in inequality 15 and this was an assumption we made in formula 17. Let us assume the opposite (inequality 15 holds), then:

$$\begin{aligned} & (d - b) \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) = 1504 \\ & (1048064) \left( \frac{32}{20000} + \sum_{i=1}^k \frac{l_i}{h} \right) = 1504 \\ & 1677 + 1048064 \sum_{i=1}^k \frac{l_i}{h} = 1504 \end{aligned}$$

The last statement is incorrect, because the left part is greater than the right part. Hence, the first assumption that inequality 15 does not hold is correct for this experiment and

$t_B = 0.73(ms)$ . Finally, time to program 1MB to the flash equals  $0.73 * 2048 = 1495.04(ms)$ , which conforms to the range given in Table II.

The results in Table II and the computations show that in case with a fast flash memory the ISP is limited by the throughput of test access path. However, the addition of the second buffer to the hybrid-online ISP mode even in case of the fast memory allows to achieve significant speed up. According to Figure 2 and 4 the possible improvement in ISP time is limited. The range of possible improvement is from 0 to  $T_1$ , where the upper limit is the time when flash is idle in the hybrid-online mode ( $T_1$ ). The time when flash is idle is called in this paper as timeout. The flash idle period in the double buffer mode is denoted here as  $T_2$ . It should be stressed that the flash is not idle during the process of copying data from buffer memory to the page buffer inside Flash. For the rest of the paper the time taken by this process is defined as  $C$ . In Figure 2 and 4 this time period is outlined as *Buffer to Page Buffer*.

Time for programming one page of data from external tester to the buffer memory (see Figure 2) equals to:  $t_B + T_1 + C$ .

The timeout in the hybrid online PCBT ( $T_1$ ) is caused solely by the time needed to transfer next page data to the buffer from external tester and this time is equal in both experiments (ISP of on-board NOR Flash using hybrid online and double-buffered hybrid online methods).

$$t_B + b \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) + C = 2.12(ms)$$

$$b \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) + C = 1.38(ms)$$

Time for programming one page of data from external tester to the buffer memory in double buffer approach (see Figure 4) equals to:  $t_B + T_2 + C = 2830/2048 = 1.38(ms)$ , where:

$$t_B + T_2 = b \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right)$$

And  $T_2 = 1.38 - t_B - C = 0.65 - C(ms)$ . In order to calculate the  $C = 0.65 - T_2$  let us subtract time to program one buffer into flash for these two approaches:

$$(t_B + T_1 + C) - (t_B + T_2 + C) = 2.12 - 1.38 = 0.74$$

$$T_1 = 0.74 + T_2 = 0.74 + (0.65 - C) = 1.39 - C$$

After substituting  $T_1$  and  $T_2$  values to  $T_1 - T_2 = 0.74(ms)$  we get:  $(1.39 - C) - (0.65 - C) = 0.74$  we get  $C = 0(ms)$ . This means that we can neglect  $C$  value in computations due to much less magnitude of measurement (not ms, but ns). It is obvious that ARM Cortex-A9 based uP is capable of copying 512 bytes of data from one memory to another in nanoseconds.

The next unresolved issue questions the possibility to overcome the bottleneck in the face of test access path. From inequality 15 it is seen that the only parameter that is not fixed in any test case is the frequency of the TCK signal. In presented experiments (Table II) the theoretical expectations

TABLE II  
ISP TIME FOR ON-BOARD FLASH DEVICES

	PCBT (ms)	2-buff. PCBT (ms)
Theoretical limit	768 - 2333	768 - 2333
Experimental data	4434	2930
Initialization	100	100

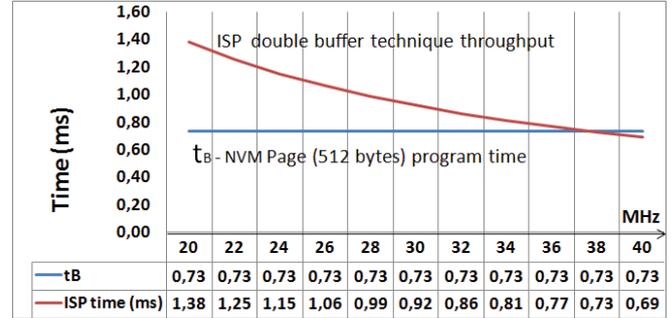


Fig. 5. ISP Time as a function of TCK frequency

are not reached with the TCK frequency fixed at 20 MHz. Figure 5 presents the chart that depicts dependency of the ISP programming time of one flash page on the TCK signal frequency.

In Figure 5 in blue is shown the  $t_B$  which is constant for the given test case and the red curve is the test access path throughput, which is a function of TCK signal frequency (horizontal axis). In other words, Figure 5 is a representation of inequality 15. At TCK frequency of 38 MHz the page transfer time becomes equal to page program time ( $t_B$ ), thus inequality 15 becomes satisfied. At higher frequency values the test access path throughput is not a bottleneck any more as the time is shorter than the  $t_B$ . Thus, execution of ISP at e.g. 40 MHz of TCK signal is excessive to achieve the shortest possible time for the given test case. Hence, the ISP time will hit the theoretically calculated limit at 38 MHz.

### C. Feasibility study

In order to see the feasibility of the hybrid online ISP solution improved by doubling the buffer we conducted experiments with various Flash devices. The results are presented in Table III. These devices are selected to study the influence of different technologies (NAND and NOR), interfaces (serial and parallel) and internal memory organisation (e.g. page size) on the overall ISP time. To make the comparison fare the test access path is similar, but not identical for all cases. The test access path cannot be identical due to the different implementation of memory controllers for NOR, SPI NOR and NAND Flash devices. However, the rest of the test access path includes uPs with identical debug ports and core architecture. These experiments are executed on the uP that implement ARMv7 [14] architecture (ARM Cortex-A9 cores) and ARM ADIV5 [15] debug interface.

The experiment with NAND Flash is conducted on the PCBA equipped with MT29T Micron [16] device and Zynq uP from Xilinx [17]. As the representative of the serial NOR

TABLE III  
ISP TIME FOR VARIOUS ON-BOARD FLASH DEVICES

	NOR [13]	qspi NOR [18]	NAND [16]
Page Size (B)	512	256	2048
$t_B$	0,73	0,50	0,39
$t_B^\dagger$	1505	2050	200
$t_{ISP_{P-OH}}^\dagger$	4434	5022	3097
$t_{ISP_{P-2BOH}}^\dagger$	2930	2973	2897

$^\dagger$  for 1 MB

Flash the device from Numonyx (N25Q128 [18]) is selected, which is paired again with Zynq uP. For the NOR Flash the data from previous experiment (see Table II) is reused.

Table III outlines the differences between various Flash devices by comparing page sizes and page programming time ( $t_B$ ). In the last two rows of Table III we present the measurement results of running hybrid online ( $t_{ISP_{P-OH}}$ ) and improved hybrid online ( $t_{ISP_{P-2BOH}}$ ) ISP solutions on the physical boards using test equipment from Goepel Electronic GMBH. In the following the measured values are used to calculate the actual time for programming a page of data, which is shown in the second row ( $t_B$ ). The computation algorithm was explained in the previous example. The actual  $t_B$  is not given in the datasheets and it varies in accordance to the number of factors (e.g. temperature, supply voltage). However, it is essential in calculations of theoretically minimal programming time under given conditions. This theoretical time value is used to assess the test access path throughput value of the ISP solution obtained in the same conditions.

Graphically the feasibility of ISP solution is presented in Figure 6. Figure 6 shows the improvement in test access path throughput that can be achieved by increasing the TCK signal frequency. In the same figure is also plotted the programming time of 1 MB data for various flash devices. When the throughput curve crosses the  $t_B$  line for 1MB it shows that the memory itself becomes a bottleneck of ISP solution. The latter shows that further increase of TCK frequency will not speed up the ISP time for particular memory. However, this figure also shows that for NAND memory the reasonable increase in TCK frequency is incapable to sufficiently increase the test access path throughput to reach the theoretical minimal ISP time.

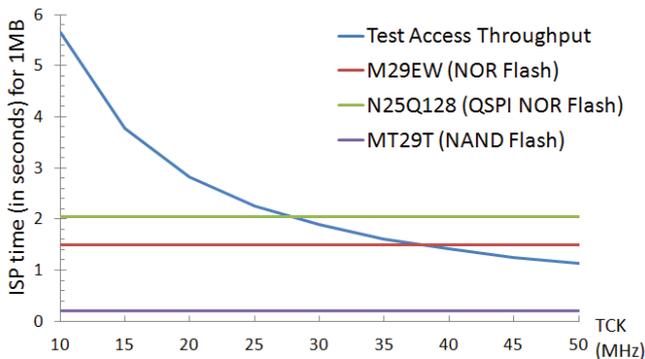


Fig. 6. Test access path throughput comparison with programming time of various Flash devices

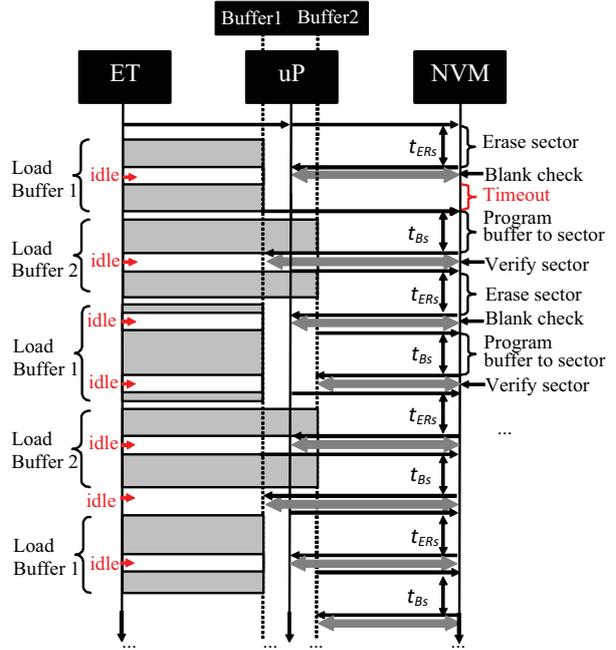


Fig. 7. Flow of ISP processes in double buffered hybrid online technique

#### D. Double Buffered Hybrid Online ISP optimizations

The experimental results in Table III show that optimization of programming step is inefficient for contemporary on-board NVM, unless the uP and PCBA design allow sufficient TCK signal frequency. Previous section studied the optimizations for programming step alone. This section considers all of ISP process steps together.

To set the goal for ISP performance we define the conditions for the ideal ISP solution. In the ideal case the NVM has to be busy through the whole ISP process. In other words the timeouts in operation of NVM are inadmissible.

Conventionally the NVM is erased, then checked for being blank, after that it is programmed and verified. In this flow the NVM is in stand by while all the data is transferred from external tester to uP and back even in case of double buffer approach.

We propose to reorganize the conventional ISP process flow as shown in Figure 7. The idea is to transfer data from external tester to buffers while the NVM is busy with erase and programming operations. The erase operation of NVM clears the content block or sector, which is more than one page. Typically, the sector erase time  $t_{ERs} = nt_{ER}$  is equivalent to the programming time for that sector  $t_{Bs} = nt_B$ , where  $n$  is the number of pages in the sector. Thus, when the ISP process is organized as shown in Figure 7 the external tester has time  $t_{ERs} + t_{Bs} = n(t_{ER} + t_B)$  to transfer data to buffer. Obviously, the size of the buffers has to be equal to the size of the sector to take advantage of the proposed ISP flow. Such big buffers may not fit into on-chip memory, so the on-board volatile memory is required, which is normally available in modern systems.

Let us discuss the flow, that is shown in Figure 7 step by step to construct the formula for ISP time calculation for the proposed solution.

- External Tester (ET) sends command to the monitor that is running inside uP to erase the sector.
- While NVM is erasing the sector, the ET starts loading the data to the buffer (*Buffer 1*).
- After the time  $t_{ER_S}$  ET pauses the data transfer and lets the monitor software (MS) inside uP to do *Blank Check* operation of the erased sector. ET is idle and listens to MS to continue data transfer for the “first” sector.
- MS reports the blank check result, and in case the sector is blank ET transfers the remaining data to the buffer. Meanwhile, the NVM is waiting for the next commands and data to program. This is theoretically the only place in the flow when NVM is in timeout.
- ET tells MS that data transfer is finished. MS programs the data page-wise to the NVM pages at operational speed of uP and NVM. This process takes time denoted as  $t_{B_S}$ .
- While the flash is busy with programming a page ( $t_B$ ), the ET transfers data to the second buffer (*Buffer 2*) for the next sector.
- After every  $t_B$  the ET halts data transfer to the second buffer and MS starts polling the NVM to check the programming operation status. After the programming of the last page in a sector is passed successfully, MS starts verification of recently programmed data against the original data that is still in the first buffer. Meanwhile, ET is polling MS for the verification process status.
- After the successful verification step MS initiates the erase operation of the next sector in NVM and reports to ET the verification result.
- ET continues the data transfer to the second buffer. ET may also start data transfer for the next (third) section if the erase operation time of the second sector is not expired.
- When  $t_{ER_S}$  expires ET halts the data transfer to allow MS do blank check of the second sector. In case the sector is blank MS starts page programming operation and ET is idling, then for time  $t_B$  ET resumes data transfer. That is again followed by programming of the next page by MS and transferring data for the following time  $t_B$  until the whole sector is programmed.
- Afterwards, MS verifies the programmed data against data in second buffer and initiates erase of the third sector. For the following time  $t_{ER_S}$  ET finishes data transfer to the first buffer and starts loading data to the second buffer.
- ... and so on, until all data is programmed

According to the performance ( $t_{ER_S}$  and  $t_{B_S}$ ) of NVM device in comparison to the test access path throughput the formula for calculation of ISP process time can take three different forms.

If inequality 18 ( $s$  - sector/block size in bytes) is satisfied, then formula 19 is justified.

$$s \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) \leq t_{ER_S} \quad (18)$$

$$t_{ISP_{2BOH}} = t_{WM} + mt_{IW} + t_{ER_S} \frac{d}{s} + t_{B_S} \frac{d}{s} \quad (19)$$

Formula 19 shows that NVM performance is a bottleneck, because programming the monitor software takes incomparably less time than the sector erase or program operations.

Second form that formula for  $t_{ISP_{2BOH}}$  can take is questioned by inequality 20. When inequality 20 is satisfied the overall ISP time is described by formula 21.

$$s \left( \frac{l}{h} + \sum_{i=1}^k \frac{l_i}{h} \right) \leq t_{ER_S} + t_{B_S} \quad (20)$$

$$t_{ISP_{2BOH}} = t_{WM} + mt_{IW} + s \frac{l}{h} + st_{IW} + t_{ER_S} \left( \frac{d}{s} - 1 \right) + t_{B_S} \frac{d}{s} \quad (21)$$

The example flow that is shown in Figure 7 reflects the formula 21. Here ISP time consists not only from monitor programming and the NVM operations, but it also takes into account the data transfer for the first sector/block.

Third form of formula for ISP time calculation is shown in formula 22. This form is considered when inequality 20 does not hold. In other words, ISP of fast NVMs, when the ET and uP is not capable of transferring data faster than NVM erases and programs it, falls into this case.

$$t_{ISP_{2BOH}} = t_{WM} + mt_{IW} + d \frac{l}{h} + dt_{IW} \quad (22)$$

#### E. Experimental results

In Table IV is presented the experimental data for full ISP process of three different on-board NVMs. The second row in this table shows the size of the sector/block of the respective NVM device. The third and fourth rows are experimentally measured data for erasing and programming sector/block. The fifth row contains the time measured for transferring data from ET to the buffer of size denoted by particular NVM sector/block size. The time measured for online hybrid ISP mode ( $t_{ISP_{OH}}$ ) is shown in the sixth row. The proposed double buffer ISP solution ( $t_{ISP_{2BOH}}$ ) results are shown in the seventh row. The last row contains calculated data for erasing and programming 1 MB of data. This data is based on the *typical*  $t_{ER}$  and  $t_B$  values from NVM device manuals.

The data presented in Table IV for selected NOR [13] device show that the proposed solution has reached the theoretical time. The difference between  $t_{ISP_{2BOH}}$  and  $t_{TYPE+P}$  is inevitable because time  $t_{TYPE+P}$  does not include the time for verification and blank check operations. By substituting values from Table IV into inequalities 18 and 20 we can conclude that in case of the given memory [13] the first form of formula for  $t_{ISP_{2BOH}}$  (formula 19) calculation should be used.

The results for qspi NOR [18] device are unexpectedly good. The measured time is even shorter than the calculated theoretical value, that is based on typical erase and program time values specified in the manual [18]. The typical time for various memory operations are specified in the manual for certain conditions (e.g. temperature). Our experiments, eventually, are executed under different conditions, that caused better performance of the particular memory. Similarly to previous case, formula 19 is proposed to be used for calculation of overall ISP time.

TABLE IV  
FULL ISP TIME\* FOR VARIOUS ON-BOARD FLASH DEVICES

	NOR [13]	qspi NOR [18]	NAND [16]
sector size	128KB	64KB	128KB
$t_{ERS}$	750	680	2
$t_{BS}$	190	120	25
$s_h^t + st_{IW}$	350	177	350
$t_{ISP_{OH}}^\dagger$	10610	14390	3120
$t_{ISP_{2BOH}}^\dagger$	7610	12880	2940
$t_{TYP_{E+P}}^\dagger$	7424	13166	169

$^\dagger$  for 1 MB at 20 MHz of TCK signal frequency

\*all time values are given in milliseconds

The ISP results ( $t_{ISP_{2BOH}}$ ) for NAND [16] memory device show that proposed optimization technique produce insufficient time improvement to compete with the performance of the selected memory. According to the results presented in Table IV for this NVM device the inequality 20 does not hold, hence, ISP time estimation is proposed to perform with the formula 22.

## VI. CONCLUSION

The contribution of this paper is twofold. The first one is the formal characterization of in-system programming of non-volatile on-chip and on-board memories. This formal basis allows to point out time losses in ISP solution and to assess the efficiency of various ISP approaches. The second contribution is the novel ISP technique that has proven to hasten programming operation to the limit that is set by the NVM performance. The proposed technique is an extension of the hybrid-online ISP mode with the double buffer concept. Even in case of exceptionally fast NVM device the proposed technique is capable to guarantee the significant speed up in comparison to the state-of-the-art solutions.

It is foreseen that the tendency of enlargement of the storage capacity of NVMs continues, thus the significance of proposed ISP technique is remarkable.

## REFERENCES

- [1] 1149.1-2001, *IEEE Standard Test Access Port and Boundary-Scan Architecture*, Std., 2001.
- [2] 1149.7-2009, *IEEE Standard for Reduced-Pin and Enhanced-Functionality Test Access Port and Boundary-Scan Architecture*, Std., 2010.
- [3] P. Maxwell, I. Hartanto, and L. Bentz, "Comparing functional and structural tests," in *Proc. of International Test Conference*, Atlantic City, NJ, USA, 2000, pp. 403 – 407.
- [4] J. Webster, B. Fenton, D. Stringer, and B. Bennetts, "On the synergy of boundary scan and emulation board test: a case study," in *Proc. of Board Test Workshop*, Charlotte, USA, 2003, p. 10.
- [5] XJTAG. (2010) High speed programming of non-volatile memories using the xjtag development system. White Paper. [Online]. Available: <http://www.xjtag.com/>
- [6] A. Tsertov, R. Ubar, A. Jutman, and S. Devadze, "Automatic soc level test path synthesis based on partial functional models," in *Proc. of 20th Asian Test Symposium (ATS)*, New Delhi, India, 2011, pp. 532 – 538.
- [7] T. Wenzel and H. Ehrenberg. (2009) Combining boundary scan and jtag emulation for advanced structural test and diagnostics. White Paper. [Online]. Available: <http://tmworld.resourcecenteronline.com>
- [8] S.Devadze, A.Jutman, A.Tsertov, M.Instenberg, and R.Ubar, "Microprocessor-based system test using debug interface," in *Proc. of 26th IEEE NORCHIP Conference*, Nov. 2008, pp. 98–101.
- [9] T.Wenzel and H.Ehrenberg. (2009) Combining boundary scan and JTAG emulation for advanced structural test and diagnostics: White paper. [Online]. Available: <http://tmworld.resourcecenteronline.com>
- [10] *XC236xE Data Sheet, V1.2 2012-06*, Infineon Technologies AG, 2012.
- [11] Keil, Tools by ARM, and ARM Ltd. (2009) Getting started. creating applications with uvision4. Manual. [Online]. Available: <http://www.keil.com/product/brochures/uv4.pdf>
- [12] Goepel Electronic Ltd. (2012) Jtag/boundary scan is probably the most ingenious test process. Web Article. [Online]. Available: <http://www.goepel.com/en/jtag-boundary-scan/boundary-scan-instruments.html>
- [13] *M29EW Datasheet, 208045-10*, Numonyx Axcell, 2010.
- [14] ARM Limited. (2011) ARM Architecture Reference Manual - ARMv7-A and ARMv7-R edition. ARM DDI 0406C, Manual.
- [15] —. (2006) ARM Debug Interface v5 Architecture Specification. ARM IHI 0031A, Manual.
- [16] Micron. (2010) M79M 2GB NAND Flash. MT29F2 Datasheet.
- [17] Xilinx. (2012) Zynq-7000 EPP Technical Reference Manual. UG585 (v1.1).
- [18] Numonyx Axcell. (2010) N25Q128 1.8V Datasheet. Rev 1.0.



**Anton Tsertov** received his M.Sc. and Ph.D. degrees in computer engineering from Tallinn University of Technology, Estonia in 2007 and 2012 respectively and currently holds the position of researcher in Tallinn University of Technology. His research interests include such topics as system and board level test, high-level system modelling, microprocessor functional and structural test.



**Sergei Devadze** has received his Ph.D. degree in computer engineering from Tallinn University of Technology, Estonia in 2009 and currently holds the position of researcher in this university. His research interests embrace such topics as usage of chip-embedded instrumentation for system test and ISP, fault tolerance and fault management architectures of digital systems, extended structural board test. He is a co-author of over 50 scientific papers in the field of digital design and test published in international journals and refereed conference proceedings.



**Artur Jutman** received his M.Sc. and Ph.D. degrees in computer engineering from TU Tallinn, Estonia in 1999 and 2004 respectively. His research interests include: embedded instrumentation for board and system test, system modeling, DFT and self-test (adding up to over 120 scientific publications). Artur Jutman has been a visiting researcher and invited lecturer in several European universities in Germany, Sweden, Poland, and Portugal. Dr. Jutman is a council member of the European Association for Education in Electrical and Information Engineering (EAEIE) and a technology development center ELIKO. He is also a member of the executive committee of the Nordic Test Forum (NTF) society. He is a managing director of Testonica Lab company, which main focus lies in the field of system test instrumentation. Dr. Jutman has been actively involved in numerous FP5, FP6, FP7 R&D projects serving as a coordinator in two of them.



**Artjom Jasnetski** received his M.Sc. degree in computer engineering from Tallinn University of Technology, Estonia in 2013 and currently he is Ph.D. student at Tallinn University of Technology. His research interests include such topics as SiP test, digital system modelling, ISP and HW driver implementation.