

Nový řídicí a dohledový systém pro experiment COMPASS

Martin Bodlák Vladimír Jarý¹ Josef Nový

¹Fakulta jaderná a fyzikálně inženýrská
ČESKÉ VYSOKÉ UČENÍ TECHNICKÉ V PRAZE
<mailto:Vladimir.Jary@cern.ch>

InstallFest 2012
Školicí centrum Silicon Hill, Praha
4. března 2012



Přehled

- 1 Systémy pro sběr dat
- 2 Sběr dat na experimentu COMPASS
- 3 Vývoj nového systému pro sběr dat
 - Vzdálené řízení
 - Nový systém pro sběr dat
 - Testy nového systému pro sběr dat



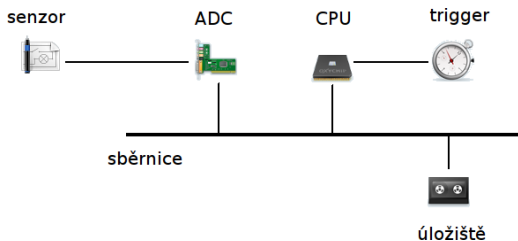
Základní pojmy

- *událost*: data popisující průlet částice systémem detektorů
- role systému pro sběr dat (*data acquisition*, DAQ):
 - 1 načtení dat z detektorů (*readout*)
 - 2 sestavení úplných událostí z fragmentů pocházejících z různých kanálů (*event building*)
 - 3 zapsání událostí do trvalého úložiště (*data logging*)
 - 4 dohled a řízení (*monitoring, run control*)
- *trigger systém*: vybírá fyzikálně zajímavé události nebo zamítá nezajímavé události
- účinnost trigger systému:
$$\epsilon = N_{\text{dobrych(vybranych)}} / N_{\text{dobrych}} < 1$$
- mrtvá doba (*deadtime*) systému:
$$D = t_{\text{system_je_vytizen}} / t_{\text{celkovy}}$$
 (je-li systém vytížen, nemůže přijímat žádné další události)



Příklad: systém s periodickým triggerem

- v podstatě se jedná o vzorkování veličiny spojité v čase
- A/D převodník digitalizuje data, CPU je načítá a ukládá
- frekvence triggeru dána dobou zpracování události:
 - je-li potřeba na zpracování 1 ms $\Rightarrow f_{trigger} \leq 1 \text{ kHz}$

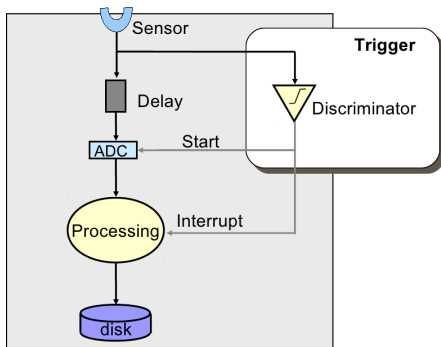


System pro sběr dat řízený periodickým triggerem



Fyzikální trigger

- data z detektorů přichází náhodně a nepředvídatelně
- potřeba mít fyzikální trigger



Systém sbírající data řízený fyzikálním triggerem podle [3]



Problémy s fyzikálním triggerem

- 1 co dělat pokud nastane nová událost a systém je zaneprázdněný:
 - přidání busy logiky: je-li busy signál aktivní, systém nepřijímá další události
- 2 jak využít neaktivní dobu:
 - vyrovnávací paměti (FIFO): vyrovnávají fluktuace na vstupu a poskytují relativně stabilní datový tok na výstupu (*derandomizace*)
- 3 jak se vypořádat s velkým množstvím kanálů ($\sim 10^6$):
 - shromažďovací moduly (např. VME desky)
 - paralelní zpracování a ukládání dat
- 4 jak zajistit minimální deadtime:
 - A/D převodník pracuje na frekvenci $\gg f$
 - zpracování a ukládání dat na frekvenci $\sim f$



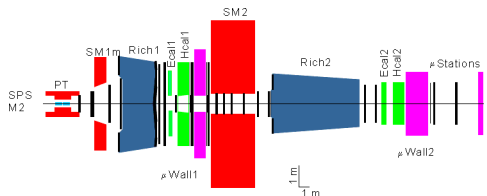
Experiment COMPASS

- *COMPASS* = COmmon Muon and Proton Apparatus for Structure and Spectroscopy
- experiment s pevným terčem na urychlovači *SPS* v laboratoři *CERN*
- vědecký program schválen v roce 1997
 - studium struktury a spektroskopie hadronů
 - experimenty s mionovým svazkem
 - experimenty s hadronovým svazkem
- sběr dat od roku 2002
- momentálně začíná 2. fáze experimentu (*COMPASS II*)
- mezinárodní projekt: 250 vědců, 29 institucí, 11 zemí



Popis experimentu

- cyklus urychlovače SPS: svazek (*beam*) není spojitý, skládá se z úseků (*spills, bursts*)
 - systém pro sběr dat používá vyrovnávací paměti pro rozložení zátěže na celý cyklus urychlovače
- interakcí svazku s terčem vznikají sekundární částice
- průlet částic detekován systémem detektorů



Systém detektorů, svazek částic dopadá na terč zleva, převzato z [4]



Vrstvy systému pro sběr dat

- 1 *primární elektronika detektorů*
 - provádí předzpracování a digitalizaci analogových dat
 - celkem zhruba 250000 kanálů
- 2 moduly *GeSiCA*, *CATCH* (VME technologie)
 - provádí načítání a shromažďování dat
 - načítání aktivováno signály z TCS (*Trigger Control System*)
 - přidání hlavičky (identifikátor triggeru, časová značka)
- 3 ROB (*readout buffer*) servery
 - slouží jako vyrovnávací paměť pro efektivní využití cyklu SPS urychlovače
 - PCI karta spillbuffer (512 MB paměti)
- 4 EVB (*event builder*) servery
 - sestavení kompletních událostí
 - zapsání souborů s událostmi na trvalé úložiště
 - uložení metadat o událostech do Oracle DB
 - doplňkové úlohy: dohled na kvalitou dat, filtr



DATE (Data Acquisition and Test Environment)

- software navržený pro experiment ALICE na LHC
- řada úprav a doplňků pro COMPASS experiment
- základní dva procesory:
 - 1 lokální shromažďovač dat: provádí načítání dat z detektoru
 - 2 globální sběrač dat: sestavuje události z fragmentů vyprodukovaných lokálními shromažďovači dat
- dobře škálovatelný a flexibilní systém:
 - režim pp (vysoká frekvence interakcí, malé události)
 - režim PbPb (nízká frekvence interakcí, velké události)
 - DAQ experimentu ALICE × malé laboratorní experimenty s jedním procesorem
- testy výkonu:
 - načítání dat: 40 GB/s
 - sestavování událostí: 2.5 GB/s
 - záznam na úložiště: 1.25 GB/s



Problémy se současným systémem

Motivace:

- časem roste frekvence trigger systému, datový tok
- 260 TB dat zaznamenáno v roce 2002, v roce 2010 již 2 PB
- vyšší datový tok \Rightarrow vyšší DAQ deadtime
- stárnoucí HW \Rightarrow vyšší poruchovost
- vývoj PCI-Express verze spillbuffer karty nákladný
- chybějící vzdálené řízení

Návrh nového systému:

- nahradit síť ROB a EVB serverů vlastním HW
- tok dat, sestavování událostí řízeno HW
- software už pouze pro řízení a dohled
- možné použití i pro další experimenty (PANDA?)



Scientific Linux CERN 5



- založen na Red Hat Enterprise Linux 5
- CERN + Fermilab → Scientific Linux
- → Scientific Linux CERN
- stránky projektu <http://linuxsoft.cern.ch/>
- RPM balíčky, balíčkovací systém yum
- AFS klient
- repozitáře s vlastním softwarem
- RSS kanály organizace (např. CERN market)



Současná řídicí místnost

Současný velín umístěn přímo v hale experimentu COMPASS:

- Výhody
 - serverovna a detektory poblíž
 - možnost přímé fyzické kontroly plynových subsystémů
- Nevýhody
 - místo, kde prochází svazek \Rightarrow problémy s radiací
 - horší dostupnost pro členy směny
 - horší ergonomie (hluk, ...)

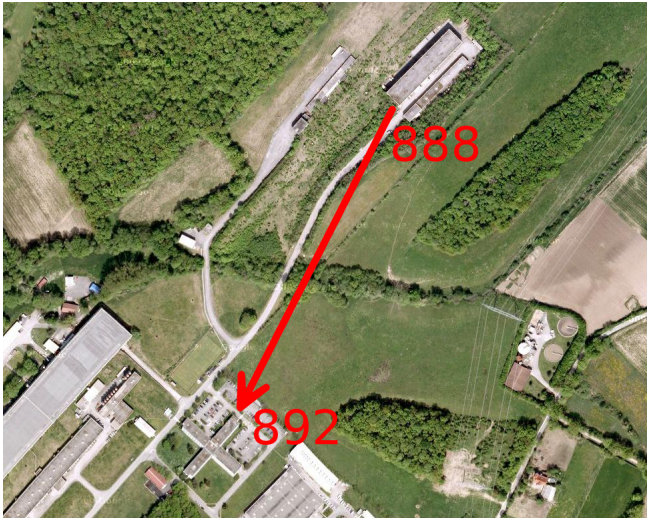
Technický koordinátor experimentu doporučil zřídit vzdálenou řídicí místnost.



Původní velín v hale experimentu



Přesun velínu



Vzdálená řídicí místnost

- vzdálený dohled a řízení experimentu COMPASS
- kancelářská budova
- vybavení:
 - 8x PC
 - 4x 24" LCD
 - 8x 22" LCD
 - IP kamery v hale experimentu
- napojeno na síť COMPASSu
- nový velín otestován
- ušetřeny finanční prostředky, které by bylo nutné investovat do přidavného stínění spektrometru
- před nasazením zbývá nainstalovat klimatizaci



Instalace nových stanic

- Scientific Linux CERN 5, 32bit
- instalační program *Anaconda*
- bezobslužná instalace pomocí *kickstart* skriptů
- různé parametry instalace podle rolí
 - řízení
 - sestavování událostí
 - souborový server
 - databázový server
 - ...
- kickstart soubory publikované v centrální databázi *AIMS*
- Průběh:
 - 1 boot po síti, stažení kickstart skriptu z databáze
 - 2 předání skriptu programu Anaconda
 - 3 načtení parametrů, pokus o bezobslužnou instalaci



Nový velín v kancelářské budově

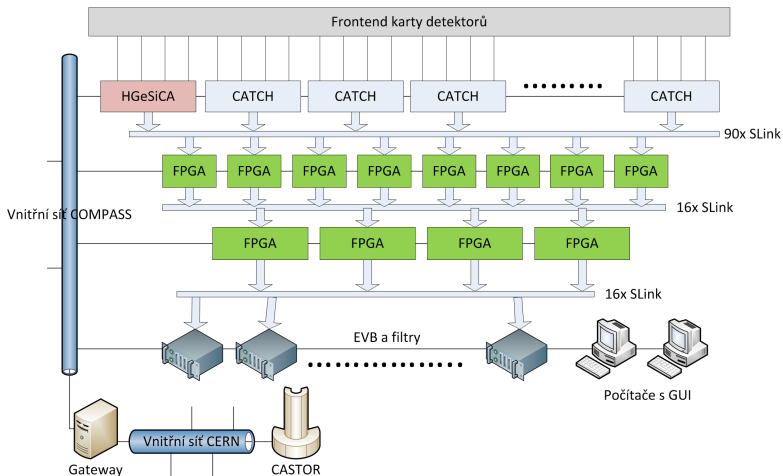


Definice požadavků na nový systém

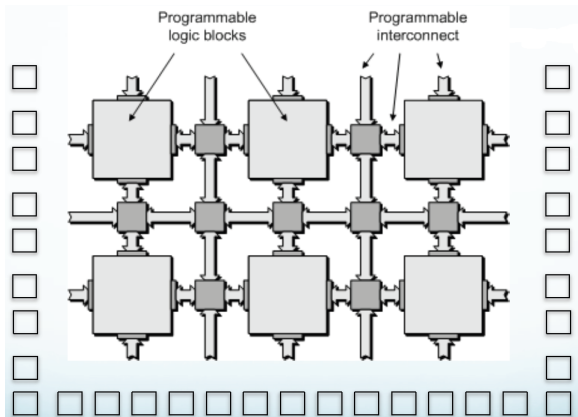
- řízení a dohled nad systémem pro sběr dat
- řízení toku dat
- jednodušší systém
- zachování stávajícího formátu dat
- použití knihovny DIM
- použití některých modulů z DATE (Murphy TV, COOOL, log book, ...)
- využití specializovaného hardwaru (FPGA karty)
- řízení v reálném čase není vyžadováno



Nová hardwarová architektura pro sběr dat



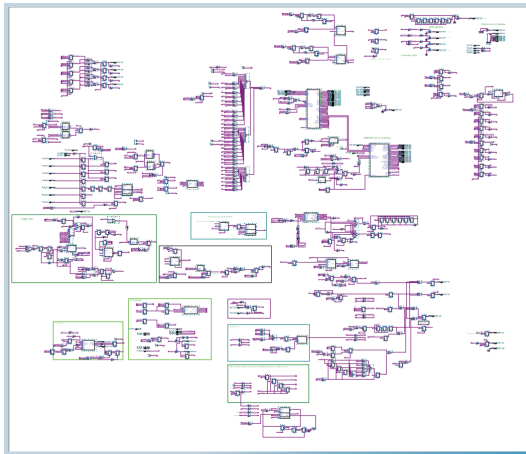
FPGA (Field-programmable gate array)



Čip programovatelný v poli (mimo továrnu), převzato z [3]



FPGA (Field-programmable gate array)



Zadání chování čipu pomocí schématu, převzato z [3]



FPGA (Field-programmable gate array)

```
architecture behavioral of VMEreg is
    signal vme_en_i : std_logic;
    signal Q : std_logic_vector(15 downto 0);

begin -- behavioral

    vme_addr_decode : process (vme_addr, vme_en) is
        variable my_addr_vec : std_logic_vector(vme_addr'high downto 0);
        variable selected : boolean;
    begin -- process vme_addr_decode
        my_addr_vec := std_logic_vector( TO_UNSIGNED ( my_vme_base_address, vme_addr'high+1 ) );
        selected := my_addr_vec(vme_addr'high downto 1) = vme_addr(vme_addr'high downto 1);
        vme_en_i <= '0' ;
        if selected then
            vme_en_i <= vme_en;
        end if;
    end process vme_addr_decode;

    reg: process (vme_clk, reset) is
    begin -- process reg
        if reset = '1' then -- asynchronous reset
            Q <= init_val;
            vme_en_out <= '0';
        elsif vme_clk'event and vme_clk = '1' then -- rising clock edge
            vme_en_out <= vme_en_i;
            if vme_en_i = '1' and vme_wr = '1' then
                Q <= vme_data;
            end if;
        end if;
    end process reg;

    data <= Q;
    vme_data_out <= Q;

end behavioral;
```

Zadání chování čipu pomocí VHDL, převzato z [3]

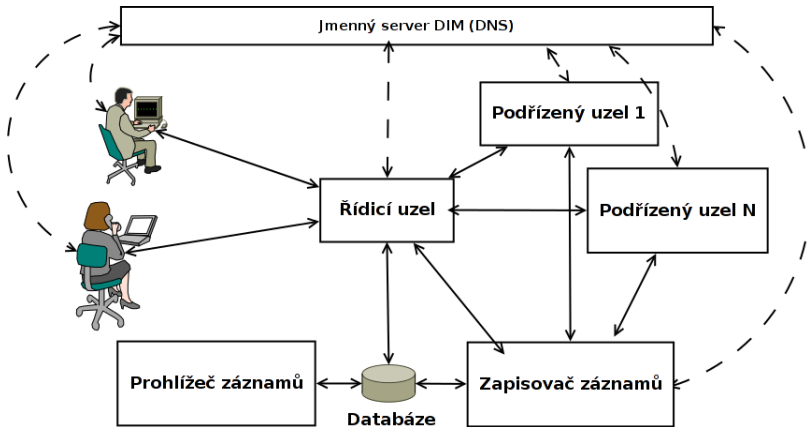


Struktura systému pro sběr dat

- Řídicí proces
 - ovládá podřízené procesy
 - komunikuje s databází
 - dostává příkazy od uživatelského rozhraní
- Podřízený proces
 - běží na specializovaném hardwaru (FPGA karta)
 - přijímá příkazy od řídicího procesu
 - poskytuje informace o stavu FPGA karty
- GUI
 - 1 řídicí rozhraní, n monitorovacích
 - přijímá informace od řídicího procesu
 - přes řídicí proces odesílá řídicí příkazy podřízeným procesům
- Message logger
- Message browser



Struktura systému pro sběr dat



Role v systému

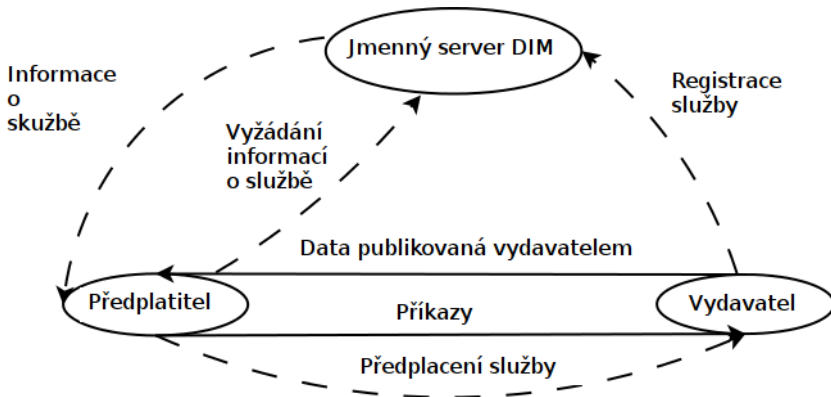


Knihovna DIM

- *Distributed Information Management*
- vývoj v CERN (původně pro experiment DELPHI)
- zajišťuje asynchronní 1 k N komunikaci po síti
- rozšíření paradigmatu klient–server o jmenný server
- postaveno na standartu TCP/IP
- rozhraní pro C/C++, Javu a Python
- multiplatformní knihovna
- používáno i v rámci DATE



Jmenný server DNS



Komunikace prostřednictvím DIM knihovny

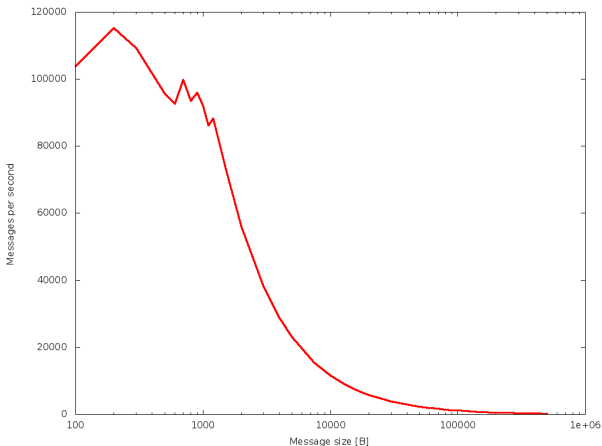


Charakteristika testovací verze

- hlavní část v QT frameworku
- nastavení a zprávy v MySQL databázi
- komunikace pomocí knihovny DIM (C++ rozhraní)
- pomocné skripty v Pythonu
- testování během zimní odstávky experimentu
- testy provedeny s parametry:
 - Gigabit Ethernet
 - 2-16 podřízených procesů na počítačích pro sestavování událostí
 - různá velikost zprávy od 100 B do 500 kB



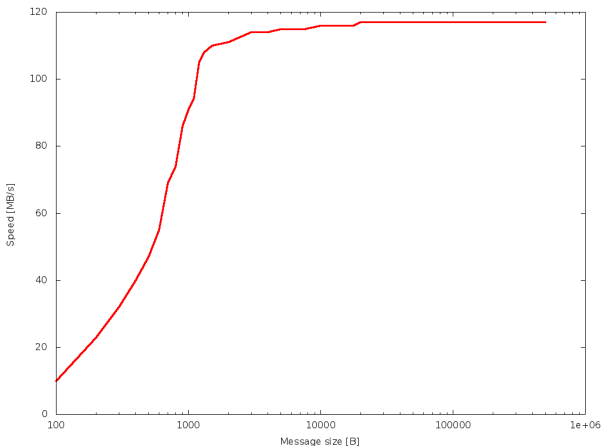
Výsledky testů (1/3)



Počet zpráv za sekundu v závislosti na velikosti zprávy



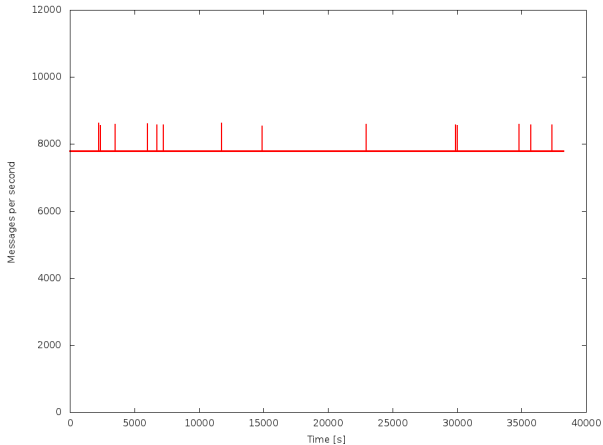
Výsledky testů (2/3)



Rychlost přenosu v závislosti na velikosti zprávy



Výsledky testů (3/3)



Test stability



Uživatelské rozhraní

The screenshot displays the COMPASS control interface with the following sections:

- Status window:** A table listing slave nodes and their status.
- Run control:** Controls for starting and stopping slaves, with a status message.
- Configuration:** Settings for the slave role, run number, number of spins, and trigger settings.
- Event size:** A large empty text area for event size configuration.
- Trigger rates:** A large empty text area for trigger rate configuration.
- Computer status:** A row of eight small monitors showing resource usage for nodes pcc0eb01 through pcc0eb08.

names	status
test001	started
test002	started
test003	started
test004	started
test005	started
test006	started
test007	started
test008	started
test009	started
test010	started

Run control: Start slaves: Connected to db successfully; Slaves started. Stop slaves: Slaves started. Buttons: Start run, Stop run.

Configuration: Role: Slaves->Master Test. Run number: 133. Number of spins: 200. Trigger settings: Random trigger. Button: Configure equipment.

Computer status:

Node	Memory	Network	CPU
pcc0eb01	2%	80%	20%
pcc0eb02	23%	90%	16%
pcc0eb03	18%	70%	17%
pcc0eb04	24%	75%	16%
pcc0eb05	19%	81%	19%
pcc0eb06	22%	77%	21%
pcc0eb07	24%	74%	18%
pcc0eb08	20%	76%	22%

Ready to start a run. Current spin: 0. Timestamp: 2010-10-20 12:06:40



Prohlížeč záznamů

Column filter: id tm dt sender severity runNum spIDNum eventNum text

tm	dt	sender	severity	runNum	spIDNum	eventNum	text
1:00 AM	11/11/11	7	FATAL ERROR	1000	1	5	Random test 1
1:00 AM	11/11/11	9	CRASH	1001	7	9	Random test 1.9
1:01 AM	11/11/11	8	ERROR	1004	13	3	Random test 2.8
1:01 AM	11/11/11	9	INFO	1004	13	12	Random test 0.9
1:01 AM	11/11/11	8	WARNING	1004	13	11	Random test 1.8
1:02 AM	11/11/11	7	INFO	1004	14	6	Random test 0.7
1:03 AM	11/11/11	8	WARNING	1004	22	9	Random test 1.8
1:03 AM	11/11/11	8	WARNING	1005	5	5	Random test 1.6
1:04 AM	11/11/11	8	FATAL ERROR	1006	14	8	Random test 3.2
1:04 AM	11/11/11	7	INFO	1009	14	7	Random test 0.7
1:04 AM	11/11/11	1	INFO	1009	14	14	Random test 0.1
1:05 AM	11/11/11	10	FATAL ERROR	1008	24	10	Random test 2.8
1:04 AM	11/11/11	8	WARNING	1005	14	22	Random test 1.8
1:05 AM	11/11/11	10	FATAL ERROR	1009	24	7	Random test 3.10
1:06 AM	11/11/11	8	FATAL ERROR	1009	25	5	Random test 1.6
1:06 AM	11/11/11	2	WARNING	1009	25	7	Random test 1.3
1:08 AM	11/11/11	7	WARNING	1009	31	10	Random test 1.7
1:08 AM	11/11/11	6	ERROR	1014	7	1	Random test 2.6
1:07 AM	11/11/11	10	WARNING	1014	12	9	Random test 1.8
1:08 AM	11/11/11	1	INFO	1014	19	10	Random test 0.1
1:09 AM	11/11/11	9	WARNING	1014	25	3	Random test 1.9
1:09 AM	11/11/11	7	ERROR	1014	25	10	Random test 2.9
1:09 AM	11/11/11	7	ERROR	1014	25	14	Random test 2.7
1:10 AM	11/11/11	4	INFO	1014	31	2	Random test 0.4
1:11 AM	11/11/11	9	ERROR	1014	37	4	Random test 1.9
1:12 AM	11/11/11	8	FATAL ERROR	1016	10	1	Random test 2.8
1:12 AM	11/11/11	7	ERROR	1016	12	10	Random test 2.7
1:12 AM	11/11/11	9	ERROR	1016	12	13	Random test 2.9
1:12 AM	11/11/11	10	WARNING	1016	17	9	Random test 1.10
1:14 AM	11/11/11	4	CRASH	1016	18	5	Random test 2.4
1:14 AM	11/11/11	2	WARNING	1016	18	8	Random test 1.2
1:15 AM	11/11/11	2	WARNING	1016	25	5	Random test 1.2
1:15 AM	11/11/11	8	FATAL ERROR	1016	24	10	Random test 1.4
1:15 AM	11/11/11	6	CRASH	1016	21	18	Random test 2.6
1:15 AM	11/11/11	5	INFO	1023	10	6	Random test 0.5
1:15 AM	11/11/11	9	INFO	1027	8	6	Random test 1.7
1:16 AM	11/11/11	9	CRASH	1027	15	6	Random test 2.9
1:16 AM	11/11/11	8	FATAL ERROR	1027	15	11	Random test 1.4
1:16 AM	11/11/11	3	INFO	1036	7	7	Random test 1.7
1:16 AM	11/11/11	3	ERROR	1036	7	11	Random test 2.3
1:16 AM	11/11/11	1	CRASH	1041	10	7	Random test 2.1
1:17 AM	11/11/11	3	FATAL ERROR	1041	10	8	Random test 3.1
1:18 AM	11/11/11	1	INFO	1041	19	5	Random test 0.1
1:18 AM	11/11/11	9	INFO	1041	19	14	Random test 0.9
1:18 AM	11/11/11	9	WARNING	1043	3	6	Random test 1.9
1:19 AM	11/11/11	2	WARNING	1043	6	5	Random test 1.2
1:19 AM	11/11/11	1	WARNING	1043	6	10	Random test 1.1
1:19 AM	11/11/11	5	INFO	1043	6	15	Random test 0.5
1:19 AM	11/11/11	8	FATAL ERROR	1043	8	21	Random test 1.8

Message filter: Info Warning Error Fatal error

Sender: test001 test002 test003 test004 test005 test006 test007 test008

Run number: Exact is 1500 From 1500 To 1500

SpID number: Exact is 55 From 25 To 55

Event number: Exact is 5 From 3 To 8

Date-time: From: Nov 11 2011 00:00 To: Nov 11 2011 00:00

Error text:







Dosažené cíle a další kroky

- 1 Analyzován současný systém pro sběr dat
 - založen na balíku DATE
 - problémy s výkonem a stabilitou
- 2 Nainstalován vzdálený velín experimentu
 - velín připraven k nasazení
 - ušetřeny finanční prostředky za přídavné stínění
- 3 Vývoj nového systému pro sběr dat
 - připraven návrh řídicího a dohledového systému
 - minimální verze tohoto systému implementována a otestována
- 4 Další kroky
 - testy na reálném HW (embedded linux na softcore procesoru)
 - rozšiřování funkcionality



Literatura

-  P. Abbon et al. (the COMPASS collaboration): *The COMPASS experiment at CERN*, In: Nucl. Instrum. Methods Phys. Res., A 577, 3 (2007) pp. 455–518
-  H. Sakulin: *Field Programmable Gate Arrays*, In: International School of Trigger and Data Acquisition, Krakow, February 2012
-  W. Vandeli: *Introduction to Data Acquisition*, In: International School of Trigger and Data Acquisition, Roma, February 2011
-  *COMPASS page* [online]. 2010.
Available at: <http://wwwcompass.cern.ch>

