



GRID технологии.  
Или как CERN си прави сметките.

Тодор Иванов

Септември 08, 2017



- HEP + IT = Проблем.



# Ускорителя LHC

- НЕР + ИТ = Проблем.

- Максимална светимост 2017:

$$1.58 * 10^{34} [cm^{-2} \cdot s^{-1}]$$

Отстояние между групите в снопа: 25 ns

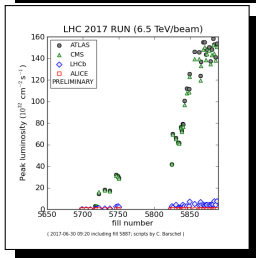
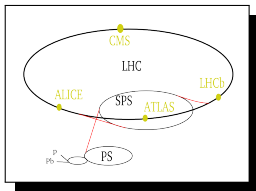
Честота на сблъсъците: 40 MHz

$$N_{events} = L\sigma_{event}$$

$$\sigma_{p-p/7TeV} \ 60mbarn$$

600 000 000 000 сблъсаци / сек

Размер на данните: 6 Gbyte/s





## CERN - Изчислителен център

*Meyrin + Vigner :*

*Memory(TiB)* 1153.72

*Numberofcores* 217530

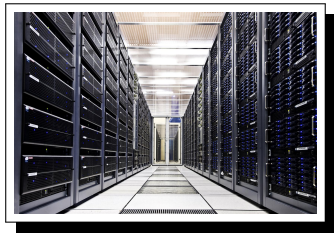
*Numberofdisks* 90179

*Numberofmemorymodules* 108743

*Numberofprocessors* 27409

*Numberofsystems* 14421

*RawHDDcapacity(TiB)* 246262

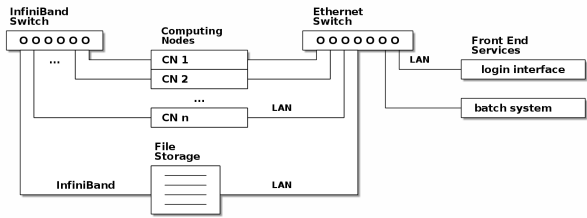


source: <http://hwcollect.cern.ch/>



# Grid: Дефиниция

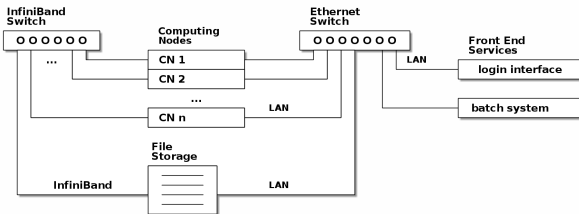
- Стандартни клъстърни системи





# Grid: Дефиниция

- Стандартни клъстърни системи

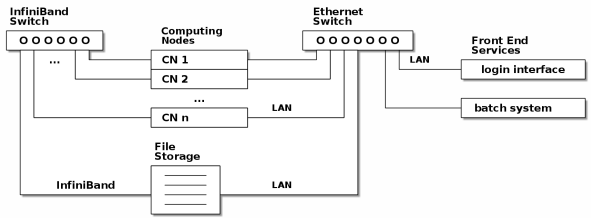


- Обединяване в мрежа от клъстърни системи:



# Grid: Дефиниция

- Стандартни клъстърни системи

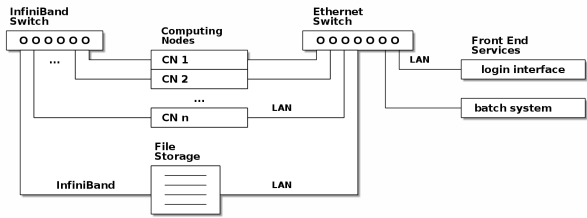


- Обединяване в мрежа от клъстърни системи:
- Нива на абстракция



# Grid: Дефиниция

- Стандартни клъстърни системи



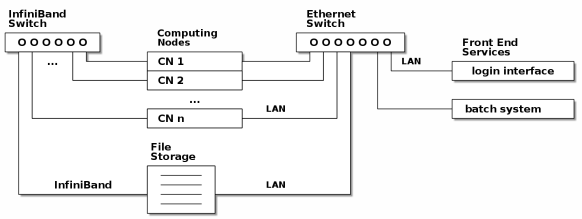
- Обединяване в мрежа от клъстърни системи:
- Нива на абстракция





# Grid: Дефиниция

- Стандартни клъстърни системи



- Обединяване в мрежа от клъстърни системи:
- Нива на абстракция





## Grid: Дефиниция

Стъпки за обединение на ниво услуги:

- Базирани на P2P комуникация - силно децентрализирани.
  - Napster - 1999г.
  - Seti@Home - 1999г. University of California, Berkeley
- Усложняване на видовете услуги и връзките между тях:



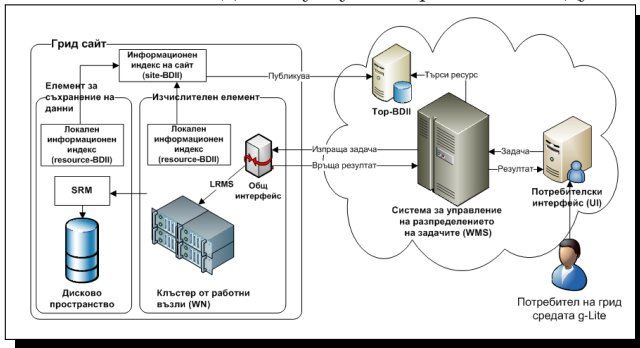
- Мидълуеър: gLite - част от проекта EGEE/EMI; Globustoolkit; VDT (Virtual Data Toolkit).



## Grid: Дефиниция

Стъпки за обединение на ниво услуги:

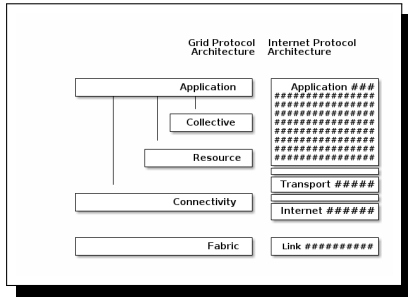
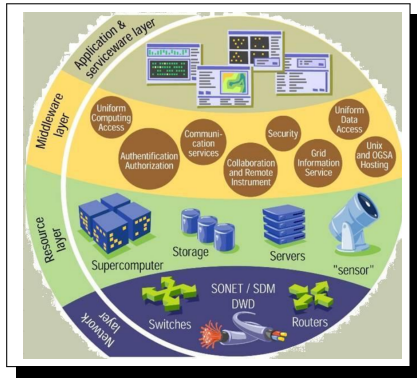
- Базирани на P2P комуникация - силно децентрализирани.
  - Napster - 1999г.
  - Seti@Home - 1999г. University of California, Berkeley
- Усложняване на видовете услуги и връзките между тях:



- Мидълуеър: gLite - част от проекта EGEE/EMI; Globustoolkit; VDT (Virtual Data Toolkit).



# Грид - Слоеве





## Grid: Дефиниция

- Виртуални организации.
- Изисквания към една мрежа от ресурси за да бъде Грид:
  - 1 Да включва в себе си и да координира споделянето на някакви географски разделени ресурси и потребители.
  - 2 Да използва стандартни, отворени протоколи и интерфейси с общо предназначение.
  - 3 Да предоставя достатъчно надежден метод за осигуряване и оценка на качеството на специализирани (QoS) ресурси.
- деф: Мрежа от ресурси с теоретично неограничен капацитет, при която основният модел на управление на ресурсите, това е принципът на “предоставяне при поискване” и притежаваща възможността за виртуализация и споделяне на самите ресурси.



## Grid: Дефиниция

- Виртуални организации.
- Изисквания към една мрежа от ресурси за да бъде Грид:
  - ① Да включва в себе си и да координира споделянето на някакви географски разделени ресурси и потребители.
  - ② Да използва стандартни, отворени протоколи и интерфейси с общо предназначение.
  - ③ Да предоставя достатъчно надежден метод за осигуряване и оценка на качеството на специализирани (QoS) ресурси.
- деф: Мрежа от ресурси с теоретично неограничен капацитет, при която основният модел на управление на ресурсите, това е принципът на “предоставяне при поискване” и притежаваща възможността за виртуализация и споделяне на самите ресурси.
- Примери:  
Particle Physics Data Grid (PPDG); EU-Datagrid ; NASA’s Information Power Grid;



## Grid: Дефиниция

- Виртуални организации.
- Изисквания към една мрежа от ресурси за да бъде Грид:
  - ① Да включва в себе си и да координира споделянето на някакви географски разделени ресурси и потребители.
  - ② Да използва стандартни, отворени протоколи и интерфейси с общо предназначение.
  - ③ Да предоставя достатъчно надежден метод за осигуряване и оценка на качеството на специализирани (QoS) ресурси.
- деф: Мрежа от ресурси с теоретично неограничен капацитет, при която основният модел на управление на ресурсите, това е принципът на “предоставяне при поискване” и притежаваща възможността за виртуализация и споделяне на самите ресурси.
- Примери:  
Particle Physics Data Grid (PPDG); EU-Datagrid ; NASA’s Information Power Grid;



## Участници в WLCG

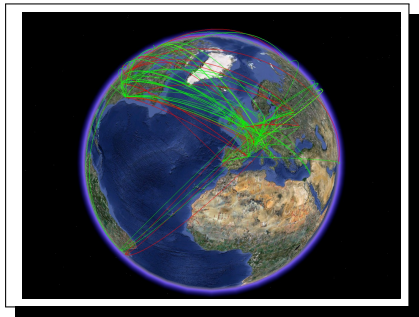


**WLCG**

Worldwide LHC Computing Grid

WLCG = EGI + OSG + NORDU GRID

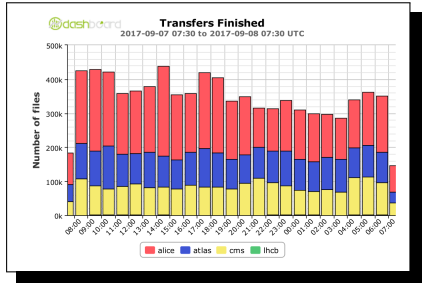
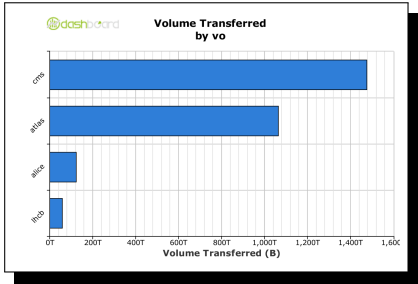
- 40 държави
- 170 изчислителни центъра
- 2 милиона изчислителни задачи/ден
- 10 Gb/s глобална скорост
- Над 10 000 потребителя





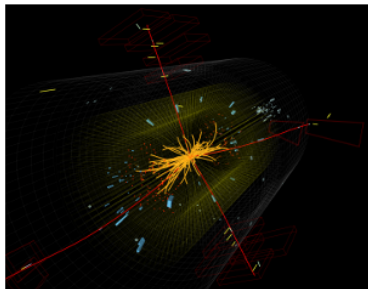


# Капацитет на WLCG



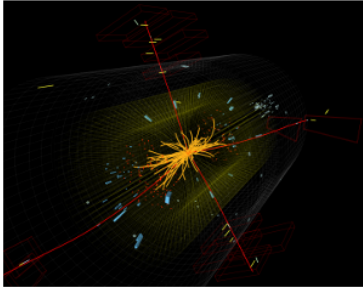


## Абстрактен модел на данните.





# Абстрактен модел на данните.

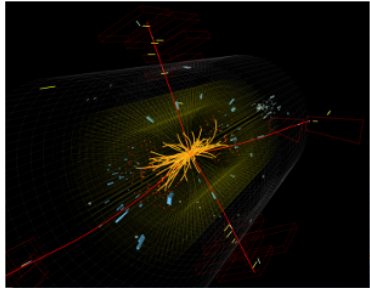


$$S[\mathbf{q}(t)] = \int_{t_1}^{t_2} L(\mathbf{q}, \dot{\mathbf{q}}, t) dt$$





# Абстрактен модел на данните.



$$S[q(t)] = \int_{t_1}^{t_2} L(q, \dot{q}, t) dt$$

```

0x01e84cf0: 0x01e8 0x87ec 0x01e8 0x85d8 0x7363 0x616e 0x0000 0x0000
0x01e84d00: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d10: 0x01e8 0x87e8 0x01e8 0x8618 0x7365 0x7400 0x0000 0x0000
0x01e84d20: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d30: 0x01e8 0x87a8 0x01e8 0x8658 0x7370 0x6c69 0x7400 0x0000
0x01e84d40: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d50: 0x01e8 0x8854 0x01e8 0x8698 0x7374 0x7269 0x6e67 0x0000
0x01e84d60: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d70: 0x01e8 0x875c 0x01e8 0x86d8 0x7375 0x6273 0x7400 0x0000
0x01e84d80: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d90: 0x01e8 0x87c0 0x01e8 0x8718 0x7377 0x6974 0x6368 0x0000

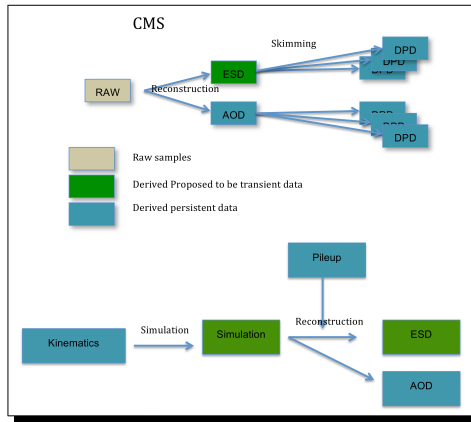
```

- Идентификатор на събитие
- Данни от детектора



# Тьер структура на данните

Тясно свързана с работните потоци.

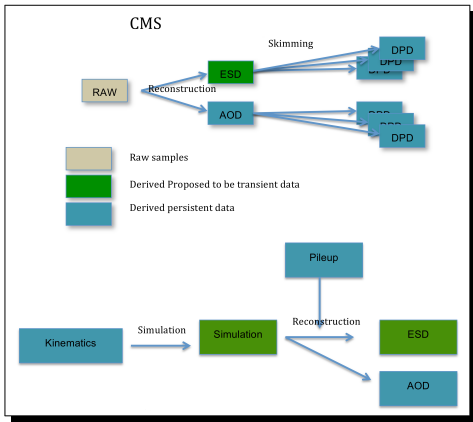


Формат	Съдържание	Размер/съб. [KB]
Типове данни от реални събития.		
DAQ-RAW	Първични - FED формат + L1 тригера	150
RAW	Първични - след онлайн форматиране - от L1 тригер и от HLT	500
ESD (RECO)	Реконструирани - трекове, вертекси, струи, частици + информацията от RAW	1000
AOD	Реконструирани - реконструираните физични обекти.	300
DPD	Извлечени Физични Данни - подобрани AOD данни, специфични за някакъв конкретен физичен анализ	tmp
NTUP	Свързани Списъци	tmp
Типове данни от симулирани събития.		
HEPMC	Изход от генератор на събития	1000
HITS	Симулирани събития - симулираната енергия, оставена в детектора от генерираните събития	tmp
RDO	Симулирани Първични данни - и (RAW-SIM-RECO)	1500



# Тиер структура на данните

Тясно свързана с работните потоци.



Формат	Съдържание	Размер/съб. [KB]
Типове данни от реални събития.		
DAQ-RAW	Първични - FED формат + L1 тригера	150
RAW	Първични - след онлайн форматиране - от L1 тригер и от HLT	500
ESD (RECO)	Реконструирани - трекове, вертекси, струи, частици + информацията от RAW	1000
AOD	Реконструирани - реконструираните физични обекти.	300
DPD	Извлечени Физични Данни - подобрани AOD данни, специфични за някакъв конкретен физичен анализ	tmp
NTUP	Свързани Списъци	tmp
Типове данни от симулирани събития.		
HEPMC	Изход от генератор на събития	1000
HITS	Симулирани събития - симулираната енергия, оставена в детектора от генерираните събития	tmp
RDO	Симулирани Първични данни - и (RAW-SIM-RECO)	1500

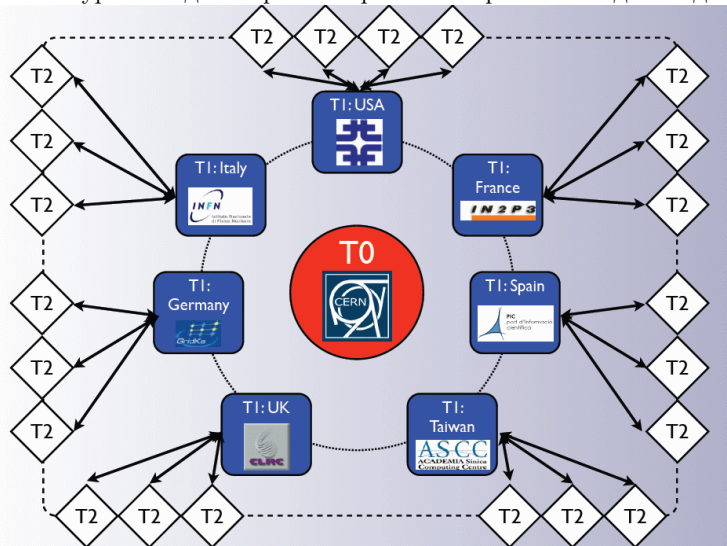
Година	Честота	Изход
2012	460 Hz	328 MB/s
2015	1000 Hz	600 MB/s

Таблица: Честота на работа на тригера от високо ниво CMS.



# Тьер структура на изчислителния модел на CMS.

Архитектурния модел е пряко свързан с абстрактния модел на данните.





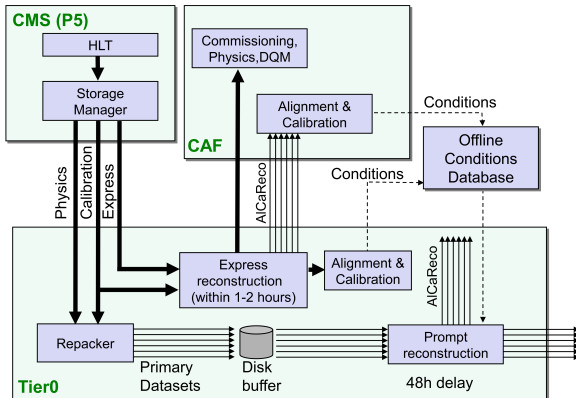
# Тьер0

## 3 потока от данни

- Физичен поток от данни с честота  $200\text{ Hz}$ .
- Експресен поток от данни с честота  $20\text{ Hz}$ .
- Поток за калибриране и наблюдение  $20\text{ Hz}$ .

## 5 работни потока

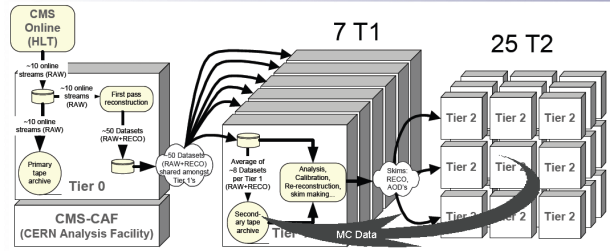
- Сливане
- Препакетиране
- Бърза реконструкция
- Експресна обработка
- Бърза калибровка





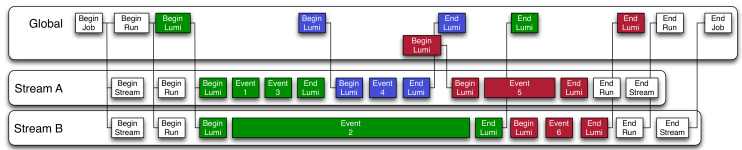


# от Тиер0 до Тиер2



Тип на изчислителните задачи:

- CPU интензивни
- Груба грануларност
- Дълго време на изчисления
- Липса на комуникация





Проблем  
Решение

Grid: Дефиниция  
WLCG  
Изчислителен модел на CMS



# Развитие на изчислителения модел на CMS

asdf



БЛАГОДАРЯ ЗА ВНИМАНИЕТО!



## Характеристиката на кълъста BG05-SUGrid

Сайт от ниво Тир3.

### 40 x COMPUTING NODES:

CPUs: 80 Intel Xeon E5345  
Cores: 160  
Threads: 320  
RAM: 640 GB  
Lan: 1 GigE

### 8 x STORAGE:

RAID arrays: 6  
Total: 72.8 TiB  
Configured Arrays: RAID6  
Distributed Filesystem: dCache

### Поддържани Виртуални Организации

- CMS - Физика на високите енергии.
- Padme - Физика на високите енергии.
- Biomed - Био-инженерни технологии, Биология, Медицина.
- Envir - Българска ВО - моделирането и защита на околната среда.
- Dteam - Изпълнява мониторираща роля.
- Ops - Изпълнява мониторираща роля.



## Производителност на BG05-SUGrid

Теста HEP-SPEC06 - оптимизиран вариант на SPEC06 (all-cp).

CPU	HEP06	Clock (MHz)	L2+L3 (KB)	Core	Memory (GB)	Mainboard
Intel Xeon E5345	58.98	2333	16384	8	16 = (8x2 FB-DDR2-677)	Supermicro X7DBE

**Таблица:** Измерване на производителността на един изчислителен възел идентичен с тези от BG05-SUGrid .

$$HEP - SPEC06(BG05 - SUGrid) = 40 * 58.98 = 2352 \quad (1)$$

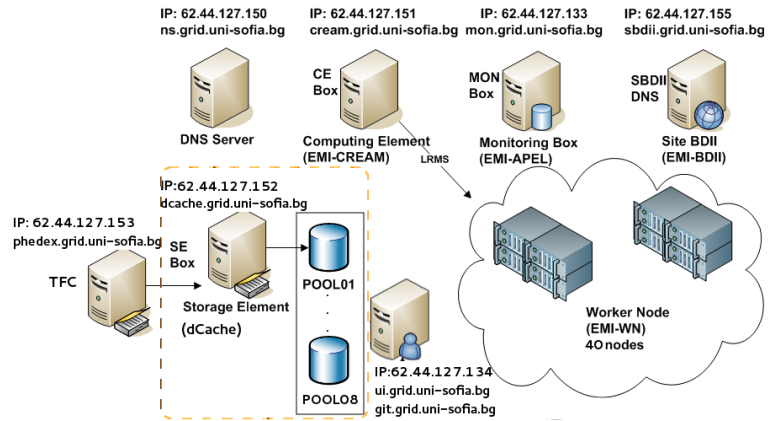
Теоретична паралелна ефективност:

$$\frac{f * N * s}{1000} = (2333 * 160 * 4) / 1000 = 1493.12 [Gflops] \quad (2)$$



# Топология на клъстър BG05-SUGrid

## GRID at University of Sofia (BG05-SUGrid)





# T3\_BG\_UNI\_SOFIA

dashb-ssb-dev.cern.ch/dashboard/request.py/siteview#currentView=default

Show 200 entries  
Copy Print Save view: default

Site Name	Visible	SAM3		Production	Analysis	Site usage		Commissioned Links Expand the details	In_rate_phedex	Out_rate_phedex	TopologyMain
		SAM3 CE	SAM3 SRM			Running	Pending				
T2_UK_SGrid_Bristol	OK	CRITICAL	WARNING	68%(293)	71%(687)	308	0	2/5 combined	39	5	2.5000
T2_UK_SGrid_RALPP	OK	DOWNTIME	DOWNTIME	84%(1085)	91%(3400)	1250	33	2/5 combined	32	4	2.5000
T2_US_Caltech	OK	OK	OK	98%(12846)	80%(11488)	3638	1745	2/5 combined	23	79	2.5000
T2_US_Florida	OK	OK	OK	94%(5676)	92%(9482)	3877	619	2/5 combined	173	123	2.5000
T2_US_MIT	OK	OK	OK	95%(7726)	85%(8348)	2798	2425	2/5 combined	14	276	2.5000
T2_US_Nebraska	OK	OK	OK	91%(6819)	81%(16815)	6820	956	2/5 combined	517	497	2.5000
T2_US_Purdue	OK	OK	OK	90%(4571)	91%(21118)	4080	726	2/5 combined	18	90	2.5000
T2_US_UCSD	OK	OK	OK	96%(5950)	75%(4259)	4432	1401	2/5 combined	59	147	2.5000
T2_US_Vanderbilt	OK	OK	OK	96%(1069)	95%(6672)	3678	1023	2/5 combined	31	201	2.5000
T2_US_Wiscnsin	OK	OK	OK	99%(8598)	84%(36516)	13873	3394	2/5 combined	32	38	2.5000
T3_AS_Panot	error	WARNING	WARNING								2.5000
T3_BG_UNI_SOFIA	OK	WARNING	OK		94%(175)	129	0		0	0	2.5000
T3_BY_NCPHEP	error	WARNING	CRITICAL								2.5000
T3_CH_CERN_Helisfebula	error	WARNING	WARNING								2.5000
T3_CH_PSI	error	WARNING	OK						0	0	2.5000
T3_CH_Volunteer	error	WARNING	WARNING		100%(1184)						2.5000
T3_CH_PKU	error	WARNING	WARNING								2.5000
T3_CO_Uniandes	error	WARNING	WARNING								2.5000
T3_ES_Oviedo	error	WARNING	WARNING								2.5000
T3_EU_Panot	error	WARNING	WARNING								2.5000
T3_FR_IPNL	OK	OK	OK		92%(1311)	430	0		1	0	2.5000
T3_GR_IASA	OK	UNKNOWN	WARNING								2.5000

Showing 1 to 150 of 150 entries First





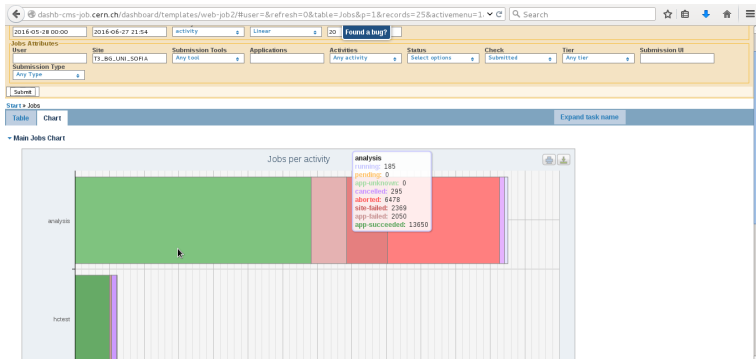


# Отработени Изчислителни задачи

The first job record is dated Tue 16 Sep 2014 12:02:47 AM EEST.  
The last job record is dated Wed 29 Jun 2016 01:46:20 PM EEST.

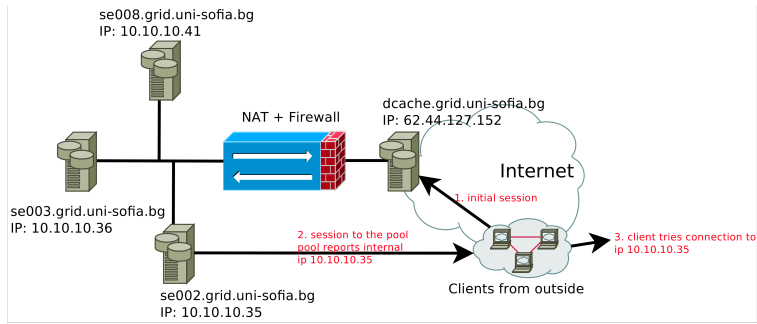
Username    Group    #jobs    days

TOTAL            -        192130    7613.08



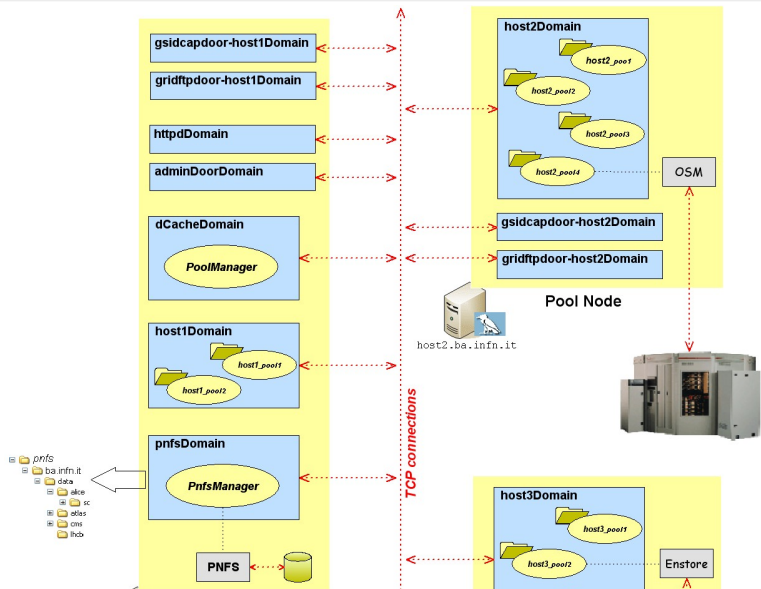


# Backup slides



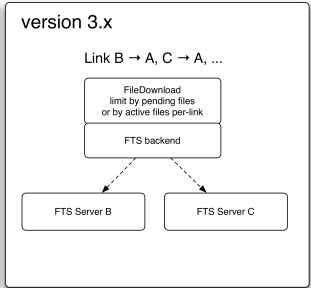
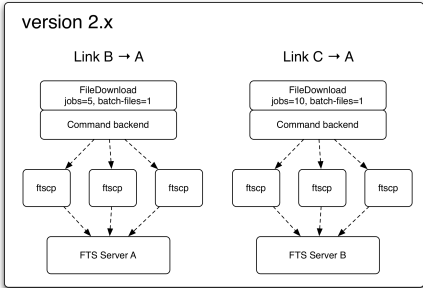


# Backup slides



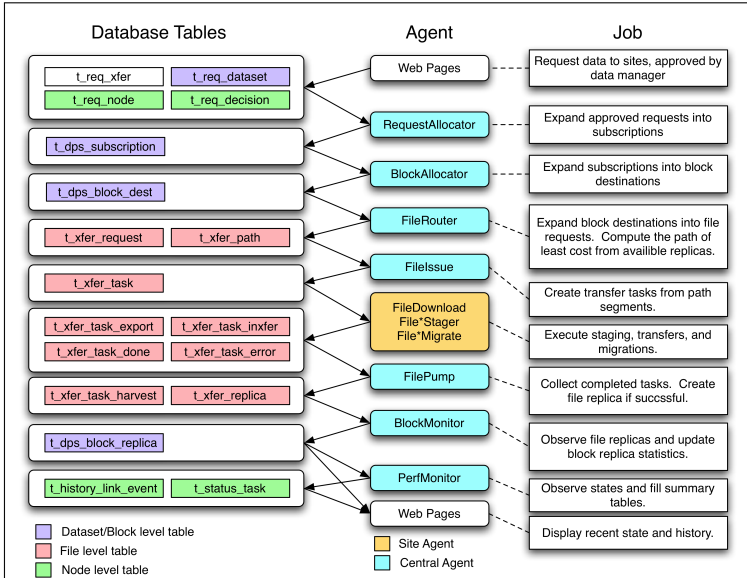


# Backup slides





# Backup slides





# Backup slides






Проблем  
Решение


Grid: Мениция  
WLCG  
Изчислителен модел на CMS



# Backup slides



## CMS Fill 5339 Report



Specific Fill: 5339  Begin: 5339 End: 5339  Stable  Last n Fills: 20  Stable  < 5338 | 5340 > [See More Options](#) [Help](#)

**CMS Fill 5339 Report**

BunchFill | LHCEvents | BeamlineLogger | ConditionBrowser | Magnet

CreateTime (localised) 2018.09.26 11:50:46  
BeginTime (stable) 2018.09.26 14:40:24  
toReady (to HV on) 0.771 minutes  
toDumpReady 3.617 minutes  
dumpReadyToDump 1.824 minutes  
EndTime (dumped) 2018.09.27 07:20:35  
Type Poses - PROTON vs PROTON  
BField 1.800 T  
Energy 6400 GeV  
InitialLumi 13062.980 /HPvsFrac4  
PeakLumi 14563.252 /HPvsFrac4  
PeakFillup (InteractionRate) 48.207  
Peak SpecificLumi 9.748 /HPvsFrac4 /HPvsFrac4  
DeliveredLumi 495.080 pb<sup>-1</sup>  
RecordedLumi 412.174 pb<sup>-1</sup>  
Efficiency by lumi 89.314 %  
Efficiency by time 95.302 %  
Physics Streams Rate 175.804 Hz  
InjectionScheme 25ns\_2250n\_2208\_2840\_2036\_2604n\_24n  
IntensityBeam1 2450.020e12 pA  
IntensityBeam2 2450.361e12 pA  
nBunchesBeam1 2220  
nBunchesBeam2 2220  
nCollidingBunches 2208  
nTargetBunches 2208  
CrossingAngle 140.045 mrad  
p 40.0 cm

Run	BeginTime	EndTime	Triggers	Lumi μs <sup>-1</sup>	Recorded μs <sup>-1</sup>	Eff %
281663	2018.09.26 14:38:23	2018.09.26 15:23:52	214090052	39.677672	31.568352	79.562
281671	2018.09.26 15:30:38	2018.09.26 15:33:41	14433	4.267520	0.002660	0.062
281674	2018.09.26 15:36:22	2018.09.26 15:54:38	704396939	18.872902	12.612572	66.845
281680	2018.09.26 16:01:34	2018.09.26 16:14:20	958800956	13.363310	8.201305	61.372
281685	2018.09.26 16:28:57	2018.09.26 16:21:42	5	2.825156	0.000000	0.000
281686	2018.09.26 16:24:04	2018.09.26 16:32:35	37278340	10.289500	5.245120	51.427
281688	2018.09.26 16:38:44	2018.09.26 16:42:55	30559983	5.969972	1.524813	25.600
281697	2018.09.26 16:47:34	2018.09.26 16:52:48	24327282	5.631297	3.576113	63.474
281699	2018.09.26 16:56:05	2018.09.27 07:11:23	3890908038	382.418148	379.311727	96.660

**Runtime chart for fill LHCFill005339**

CMS: Fill 5339 Luminosity

CMS: Fill 5339 Instantaneous Luminosity ■ CMS Online Lumi

CMS: Fill 5339 Lumi per Crossing

— Lumi per bx

— Soec Lumi per bx



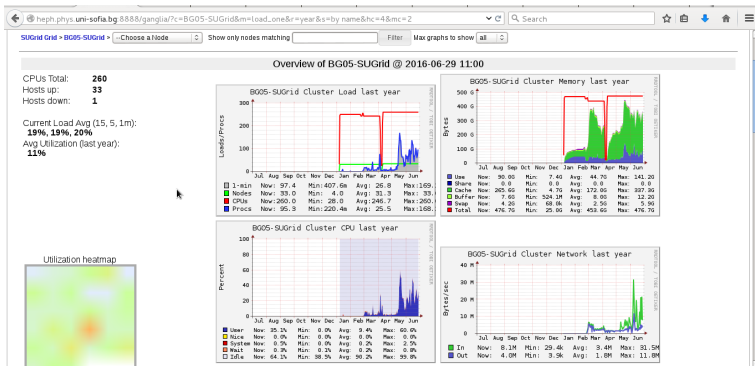
# Backup slides







# Backup slides





# Backup slides

accounting.egi.eu/egi.php?ExecutingSite=BG05-SUGrid&query=sumlap&startYear=2015&startMonth=1&endYear=2015

EGI ACCOUNTING PORTAL

CSGA

GENERAL Menu | VO MANAGER Menu | VO MEMBER Menu | SITE ADMIN Menu | USER Menu | SERVICES | METRICS TO PDL | LOGS

The following table shows the distribution of Total elapsed time grouped by SITE and VO (only information about TOP 10 - ordered by CPU time- VOs (Excluded dman and ops VOs) is showed in detail. The rest of VOs will be grouped in a new category).

Executing SITE	Elapsed	CPU	IO	Network	Excluded dman and ops VOs
BG05-SUGrid	29,280	132,312	4	27,849	38.89%
Other	19,240	307,210	4	27,849	
Percentage	11,236	18,726	8,009		

Click here to a table group of this table  
Click here to a table group of this table

Chart showing the Cumulative Total elapsed time grouped by SITE and VO (only information about TOP 10 - ordered by CPU time- VOs (Excluded dman and ops VOs) is showed in detail. The rest of VOs will be grouped in a new category).

Developed by CSGA EGI Year 1 usage | 2015-2016 | CPU VO | http://egi.eu/ | 2016-09-20 14:04

Legend: BG05-SUGrid Cumulative Total elapsed time by SITE and VO (Excluded dman and ops VOs)

Developed by CSGA

CSGA



# Backup slides

accounting.egi.eu/repintngi.php?query=sumcpu&startYear=2015&startMonth=7&endYear=2016&endMonth=6

EGE ACCOUNTING PORTAL

GLOBAL View VO MANAGER View VO MEMBER View SITE ADMIN View USER View REPORTS METRICS PORTAL LINKS

The following table shows the relative consumption of resources between NGIs. The number in bold shows the consumption of resources owned by the NGI corresponding to that column consumed by that row. Hovering over a cell shows the consumer and provider NGIs for that cell. The blue (or column) percentage represents the percentage over the total resources consumed that are provided by a NGI. The red (or row) percentage represents the percentage of resources var the total consumed by a NGI.

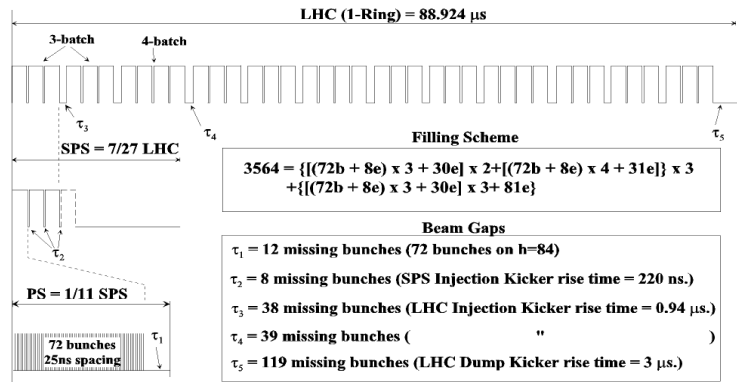
Used by	AfricaAraba	AsiaPacific	CERN	GLN	ILDIR	INDG	NGI_AEGIS	NGI_ARMGWB	NGI_BG	NGI_BY	NGI_CH	NGI_CHINA	NGI_CZ	NGI_DE	NGI_ET	NGI_FRANCE	NGI_GE	NGI				
AfricaAraba									1311		158							4272	29	1256		
AsiaPacific		2693308 20.5% / 80.48%							8.4% / 0.31%		8%	0%	0%	0.92%				0%	8.88%	0%	0.04%	
CERN	535120 4.08% / 0.83%	48734674 34.36% / 2.64%	369025948 28%				46525	261897	84.88%	0.02%	32976416 99.99%	7533878 1.75%	3868969 0.4%	211583919 51.86%	1.02%	88.51%	11.28%	222579193 79.31%	11.83%	16425	47.91%	
GLN																						
ILDIR																						
INDG																						
NGI_ARMGWB							35188 0.33%	100%														
NGI_BG								91 0.08%	100%													
NGI_BY																						
NGI_CH											1573 0%	87.88%										
NGI_CHINA																						
NGI_CZ													22724431 5.54%	99.99%								
NGI_DE		28924 0.23%	3.08%								10866 0.02%	0.01%		2424867 0.6%	0.02%	81.62%		14290 0.02%	0.06%	141	0%	
NGI_ET																						
NGI_FRANCE		7646 0.02%	0.02%					14886 0.08%	0.02%		45724 0.0%	0.08%	2449435 4.49%	2837229 1.39%	3.44%			5523852 20.86%	0.88%	4910 0.03%	28844 8.08%	
NGI_GE																				9 0.07%	100%	
NGI_GRIDNET	28 0%	357 0%	0.02%					21 0.01%	0%	5697 0.02%	0.01%		26 0%	972 0%	0.1%		337 0.01%	0.02%	441 0%	0.04%	3 0.07%	97692 28.08%

CERN consumes 94.06% of the resources provided on NGI.BG



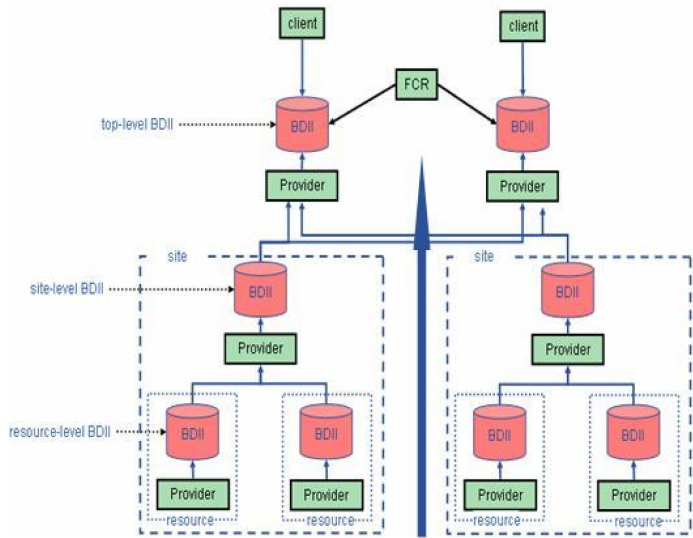
# Backup slides

## Bunch Disposition in the LHC, SPS and PS



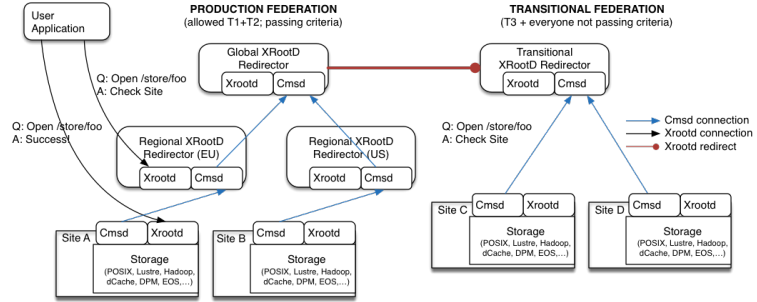


# Backup slides





# Backup slides



**Xrootd Global Redirector**

