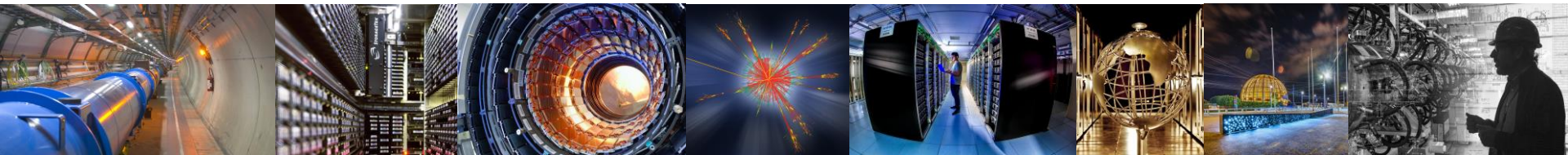


# ИТ в ЦЕРН

## Разпределени изчисления в ATLAS

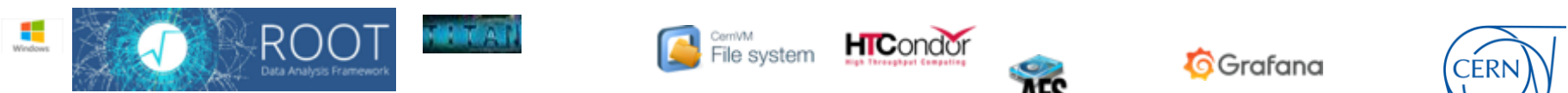


*Др. Иван Глушков  
Тексаски университет / АТЛАС  
Българска учителска програма  
ЦЕРН, юли 2019*

# ИТ в ЦЕРН

- Обслужва ~15000 потребители в цял свят
  - Физици, инженери, програмисти, администрация, финансисти
  - Поддържа пълен набор от ИТ решения за съответните задачи
  - (когато е необходимо) Разработва специфични за ЦЕРН решения
  - Оптимална цена
  - Изследване на нови насоки в ИТ





# Технологии в ЦЕРН..



Разпределени изчисления в ATLAS. БУП 2019

... за които ще си говорим само ако вие искате

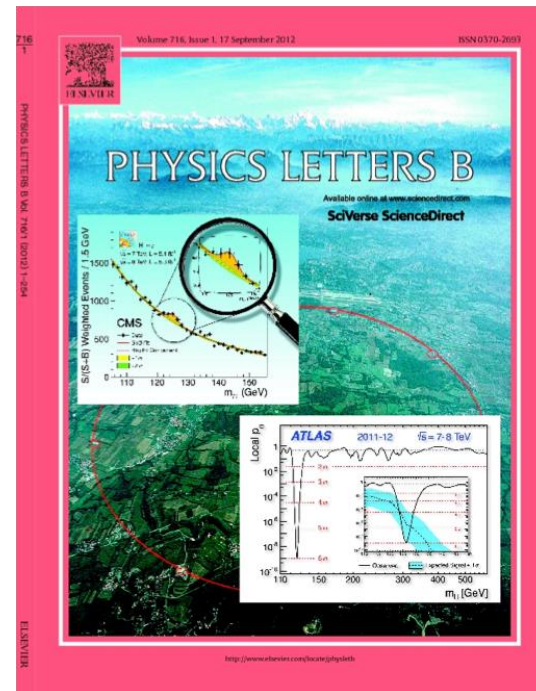
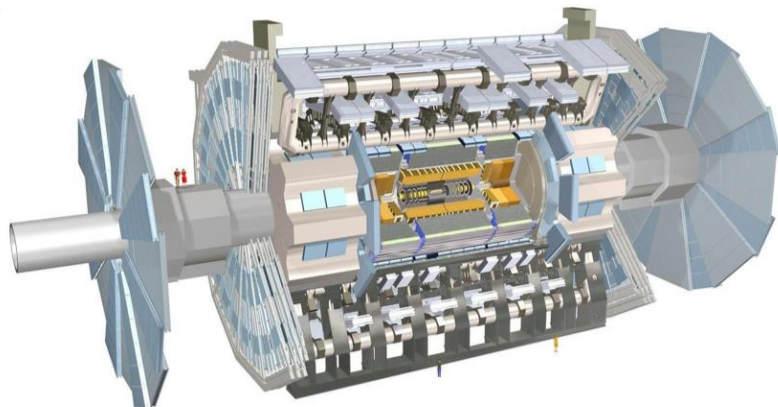


# Разпределени изчисления (в АТЛАС)

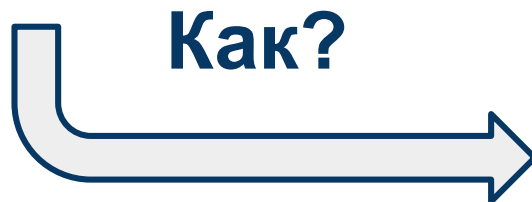
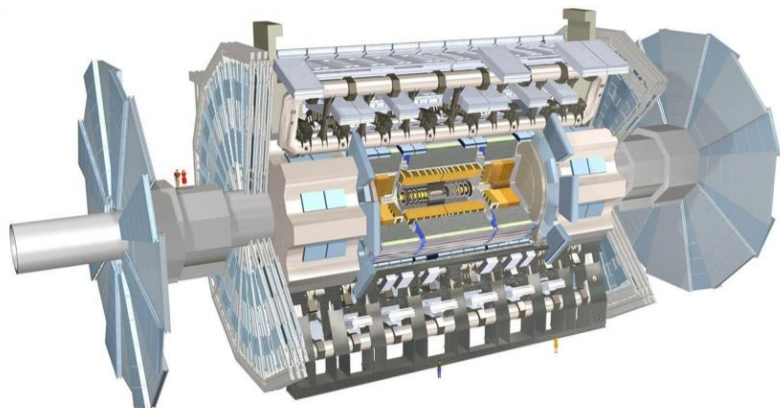


**Какъв е проблема?**

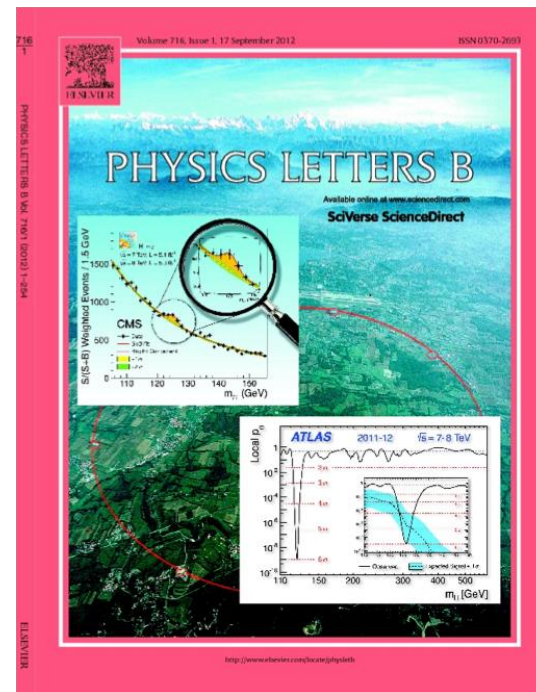
# От ATLAS до научна публикация



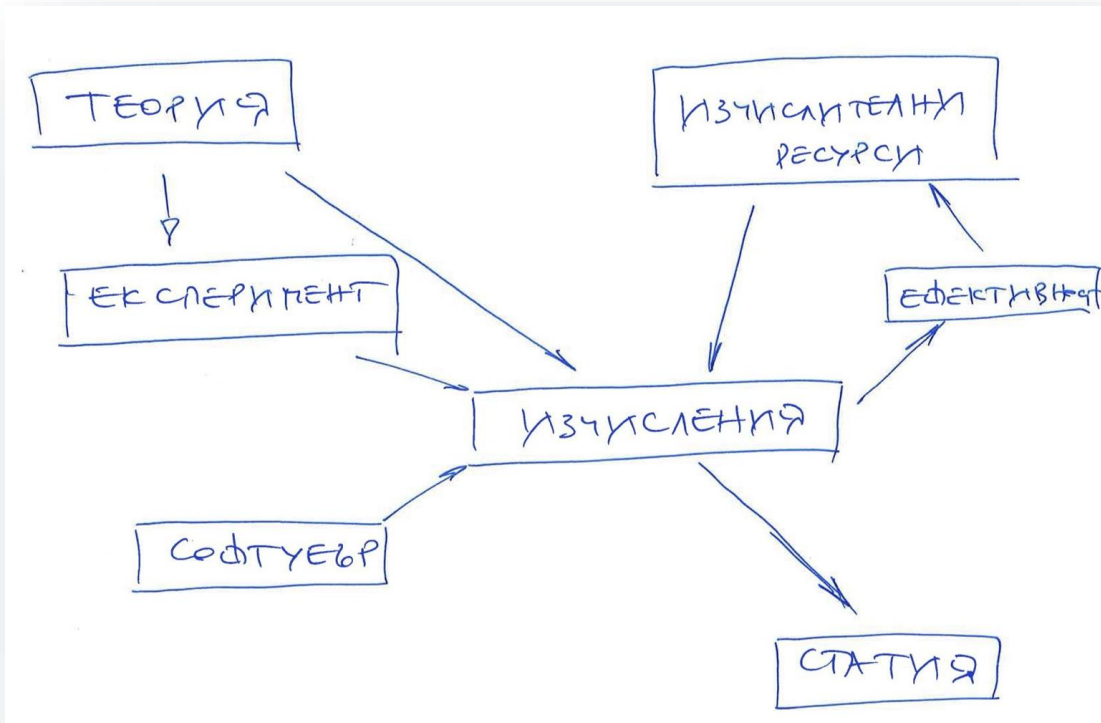
# От ATLAS до научна публикация



.. и бързо моля



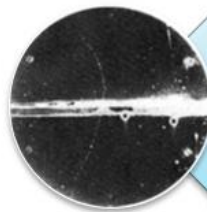
# Парадигма на науката





# Какво е нужно за откритие?

(обработка на данни)



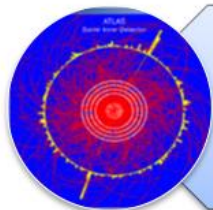
През 1930-те

- ~2 учени от една държава
- Лист и химикал



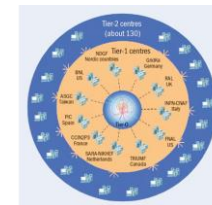
През 1970-те

- ~200 учени ~10 държави
- Мейнфрейм



Днес

- ~3000 учени ~100 държави
- **Разпределени изчисления**

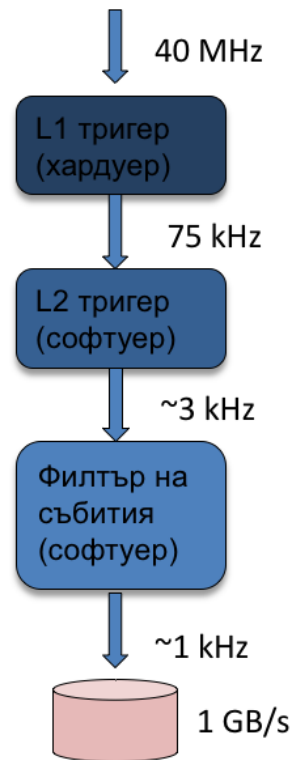
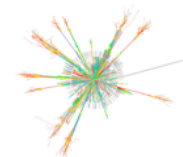




# Данните

# Колко са данните от ATLAS

- Произвеждаме: 40 MHz x 1.5 MB = 60 TB/s
- Данни == C++ обекти
- Време за обработка: ~дни т.е. ~ O(MB - GB)
- Пресяване в реално време:
  - L1 тригер - на експеримента
  - L2 тригер - ~120 000 ядра



```
SELECT SUM(bytes)/1000/1000/1000/1000/1000 FROM atlas_rucio.dids  
WHERE did_type='D'  
and datatype = 'RAW'
```

Query Result x

SQL | All Rows Fetched: 1 in 68,69 seconds

SUM(BYTES)/1000/1000/1000/1000/1000
58.43299170582995

## 58 PB!

# КАК?!

## 58 PB!?

Колко време?  
Ами ако сбъркам?  
Ами ако колегата  
направи по добър  
алгоритъм за  
мюони?  
Всичко отначало?!



# Колаборация и оптимизация

**58 PB!?**

Колко време?  
Ами ако съберам?  
Ами ако колегата  
направи по добър  
алгоритъм за  
мюони?  
Всичко отначало?!

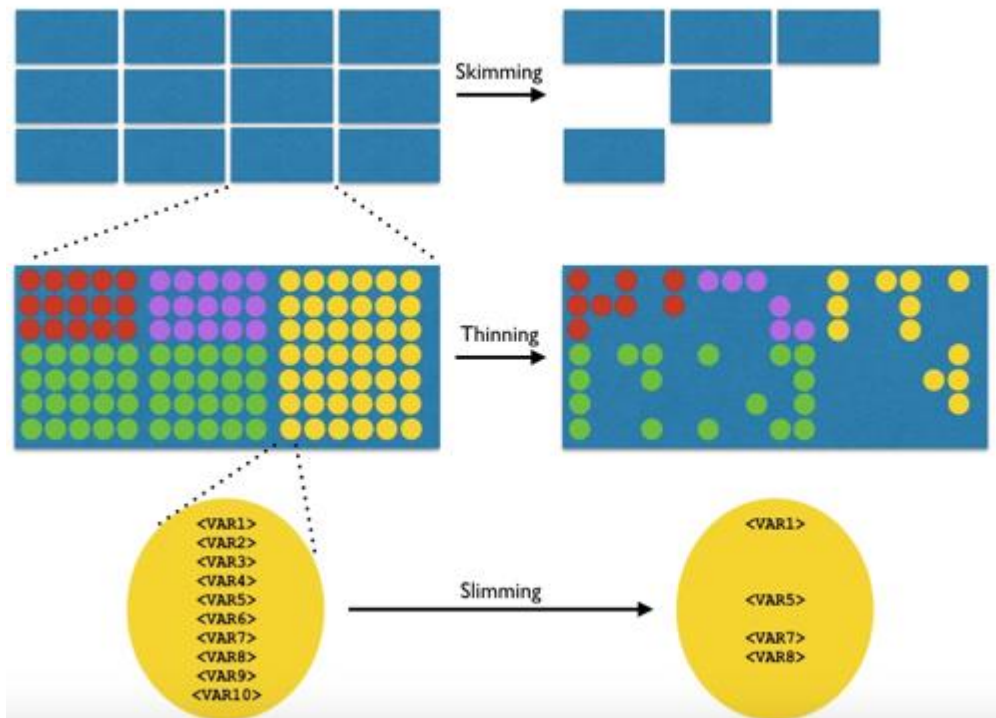


- Групиране на анализите по теми - Хигс, Суперсиметрии, Екзотики
- Селектиране само на данните които са релевантни за дадения анализ

**GB / MB**

# Селектиране на данни

## (ATLAS Derivation Framework)



## Премахване на ненужната информация

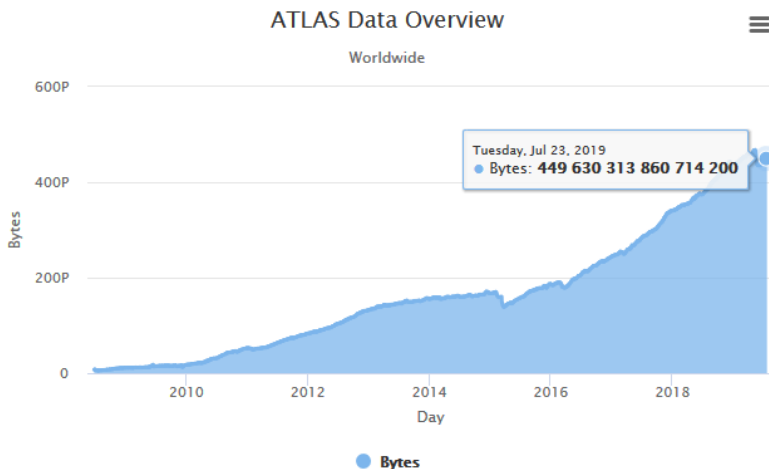
- Цели събития
- Части от събитията
- Характеристики на събитията
- За всеки анализ - отделни данни

Графика: James Catmore

# Истина и теория

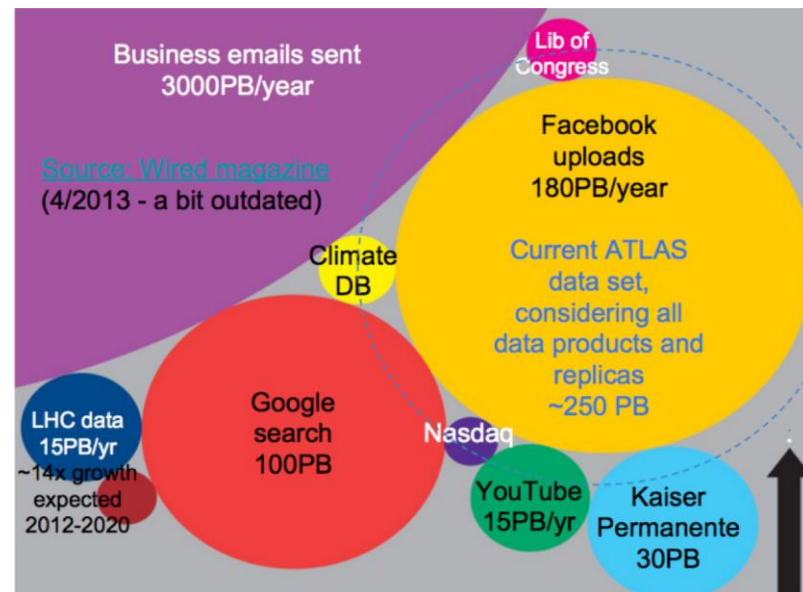


# Всичко общо..



## 450 PB!

(420 PB преди 6 месеца)







# Изчисления

(колко и какви)

**Колко изчислителни  
ресурси ни трябвават?**

**“Колкото повече, толкова  
повече!”**

**Мечо Пух, 1966**



**“Колкото повече, толкова  
повече!”**

**Мечо Пух, 1966**



**“Има само едно нещо, което е по-хубаво от гърненце  
мед.. И това са две гърненца мед.”**

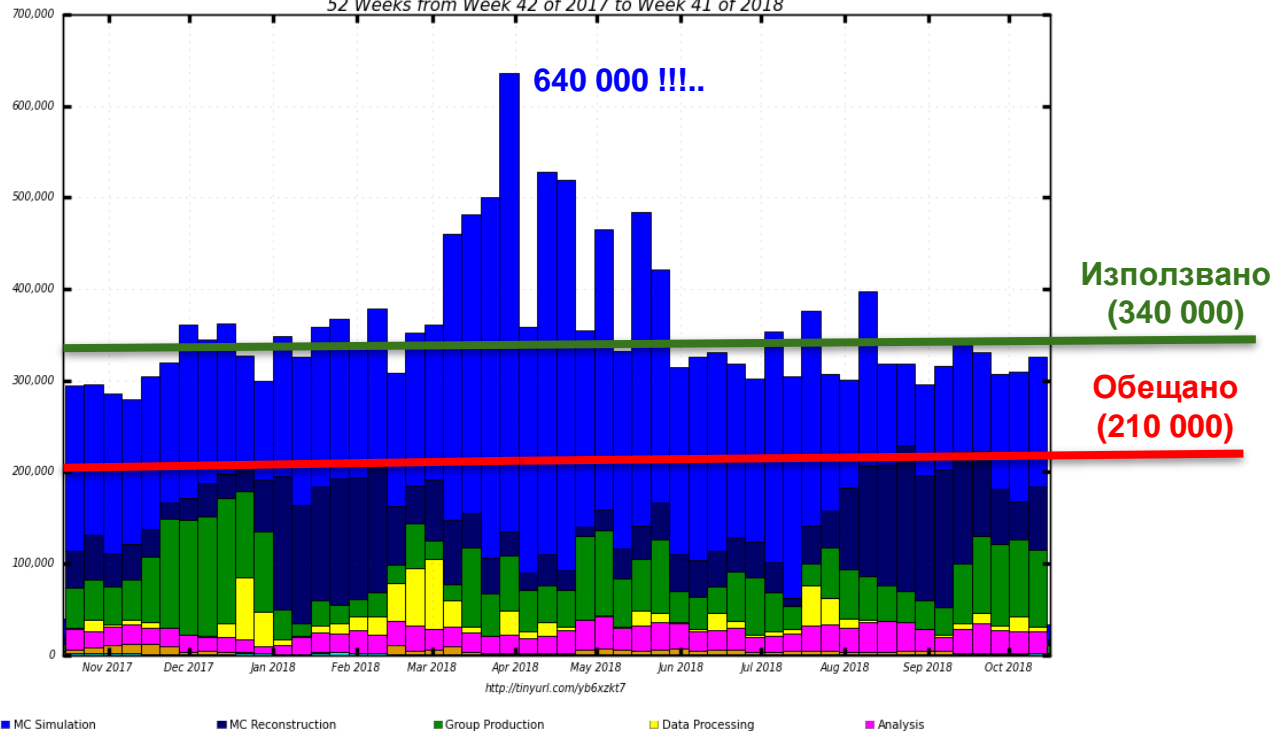
**Мечо Пух, 1966**

# Колко използваме?



## Slots of Running Jobs

52 Weeks from Week 42 of 2017 to Week 41 of 2018



Maximum: 636,447, Minimum: 0.00, Average: 340,411, Current: 33,035

## Речник

- Слот - място за един процес заделено от компютърна ферма
- Задача (Job) - количеството “работа” което трябва да се сметне на един слот

## Колкото повече, толкова...

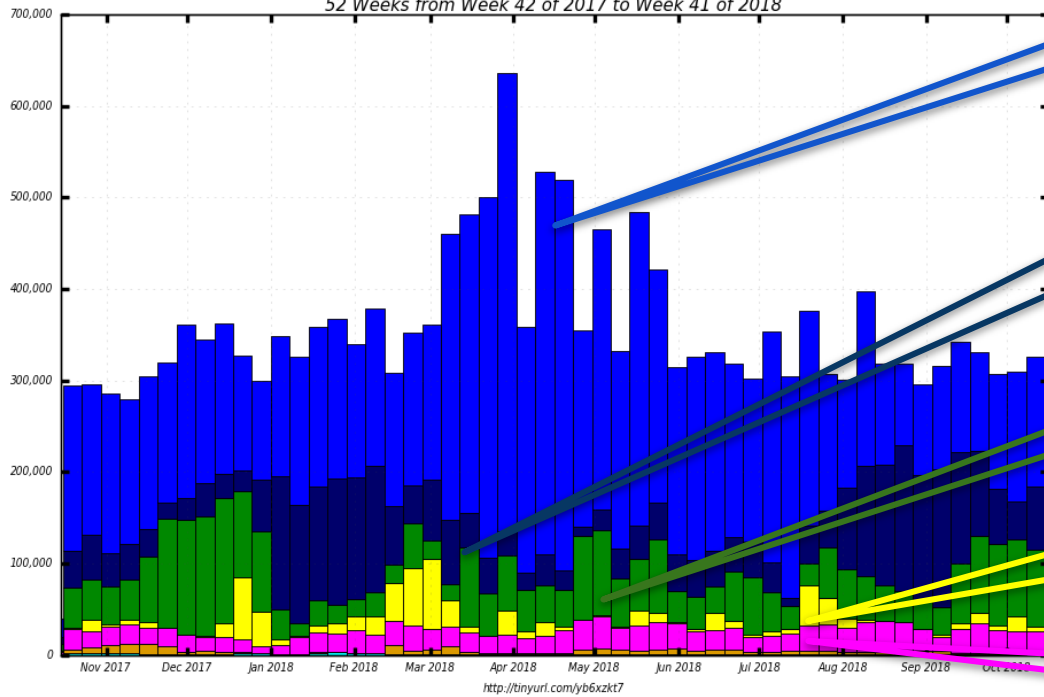
- По-бързо
- Изследване на повече нови опции - генератори, конфигурации...

# За какво ги използваме?

dashboard

Slots of Running Jobs

52 Weeks from Week 42 of 2017 to Week 41 of 2018



<http://tinyurl.com/yb6xzk7>

■ MC Simulation    ■ MC Reconstruction    ■ Group Production    ■ Data Processing    ■ Analysis  
 ■ TO Processing    ■ Others

Maximum: 636,447, Minimum: 0.00, Average: 340,411, Current: 33,035

Детекторна  
Симулация

Реконструкция

Пресяване

Преработка

Анализ



# Изчисления

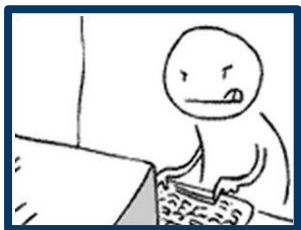
(къде)

# Разпрпределени изчисления в ATLAS: Обзор

- **Експерти: > 200 (в различна степен на заетост)**
  - Системата не може да работи без тях
  - Разработка и интеграция на нововъведения
  - Ескплоатиране / опериране на системата
- **(Централизирани) услуги**
  - Разпределяне на задачи - PanDA (ATLAS)
  - Разпределяне на данни - Rucio (ATLAS)
  - Прехвърляне на файлове - FTS (WLCG)
  - Разпределяне на софтуер - CVMFS (WLCG)
- **Отдалечени компютърни центрове (сайтове)**
  - GRID
  - Облаци, суперкомпютри, доброволчески изчисления
  - Хардуер: ресурси за обработка и съхранение на данни
  - .... и експерти на местно ниво!



# Потребител и разпределени изчисления: Принцип на работа



Къде / кои са данните които  
искам да обработвам?

Система за  
разпределяне на данни  
(Rucio)

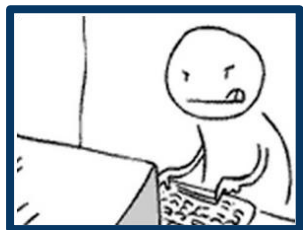
Данни X, които се намират  
на сайт Y

## Речник

- Сайт - (отдалечен)  
компютърен център

# Потребител и разпределени изчисления:

## Принцип на работа



Къде / кои са данните които  
искам да обработвам?

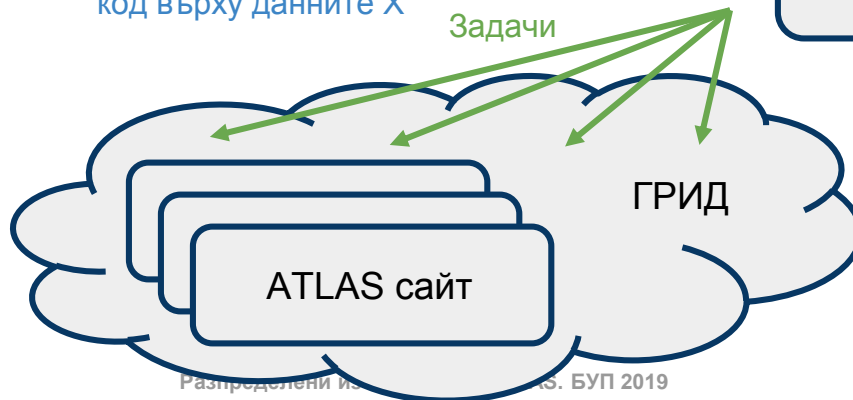
Система за  
разпределяне на данни  
(Rucio)

Данни X, които се намират  
на сайт Y

Искам да изпълниш моя  
код върху данните X

Система за  
разпределяне на задачи  
(PanDA)

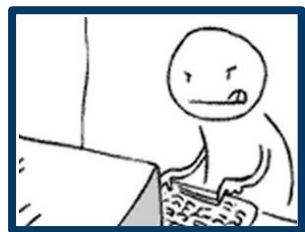
Задачи



### Речник

- Сайт - (отдалечен)  
компютърен център

# Потребител и разпределени изчисления: Принцип на работа



Къде / кои са данните които  
искам да обработвам?

Система за  
разпределяне на данни  
(Rucio)

Данни X, които се намират  
на сайт Y

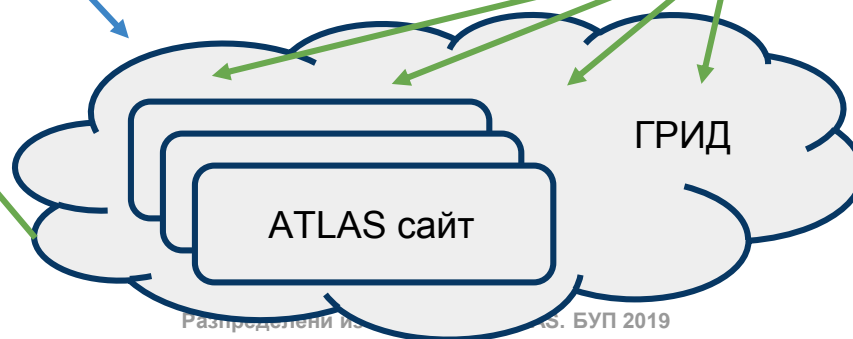
Система за  
разпределяне на задачи  
(PanDA)

Искам да изпълниш моя  
код върху данните X

Задачи

Резултати?

Резултати / Още  
не е готово



## Речник

- Сайт - (отдалечен)  
компютърен център

# Какво е GRID



GRID е **технология** която позволява оптимизиран, оторизиран достъп до ресурси - компютърни и за съхранение на данни - принадлежащи на различни собственици.

# Какво е GRID II

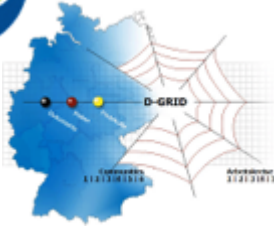


- Идеята се ражда през 90-те
- Аналог на мрежата за електроснабдяване

# Какво е GRID III

- Софтуер или “GRID middleware” - дава достъп до хетерогенните ресурси на GRID с общ интерфейс
  - CE (Computing Element) - достъп до изчислителни ресурси
  - SE (Storage Element) - достъп до ресурси за съхранение на данни
- Потребителски достъп - идентификация и упълномощаване
  - Няма как да дадем логин и парола на x10000 хора на x100 сайта
  - Дигитални сертификати (x509-базирани)
    - Раздавани от сертифицирани “авторитетни източници” (certification authority)
    - Важи за всеки GRID сайт
    - Локално на сайта всеки сайт съответства на локален акаунт (и дефинира правата)

# GRID-ove, GRID-ove..

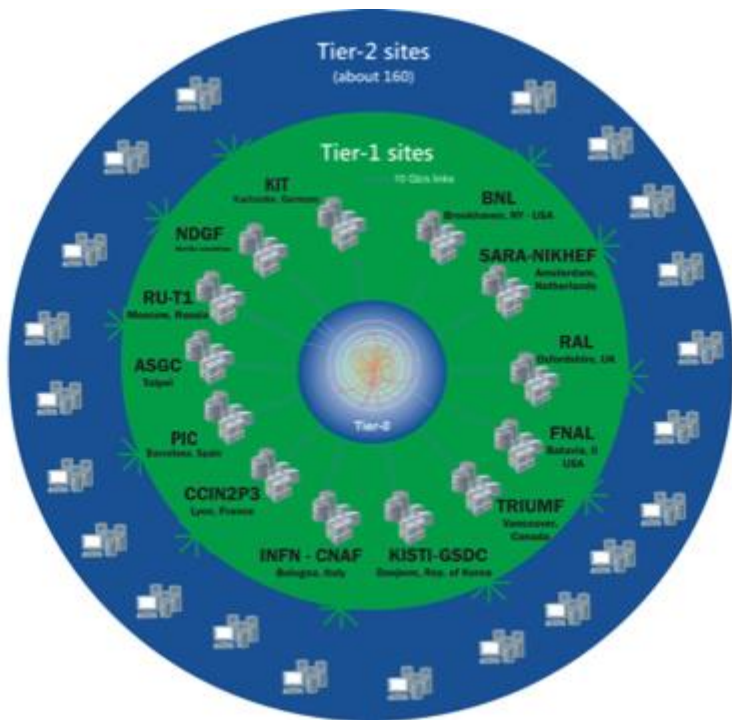


KMI

GFK



# WLCG (The Worldwide Computing GRID)



Участващи GRID-ове:

- EGI (European GRID Infrastructure) - Европейска GRID инфраструктура
- OSG (Open Science Grid) - GRID за общодостъпна наука (САЩ)
- NDGF (Nordic Data Grid Facility) - Скандинавски GRID



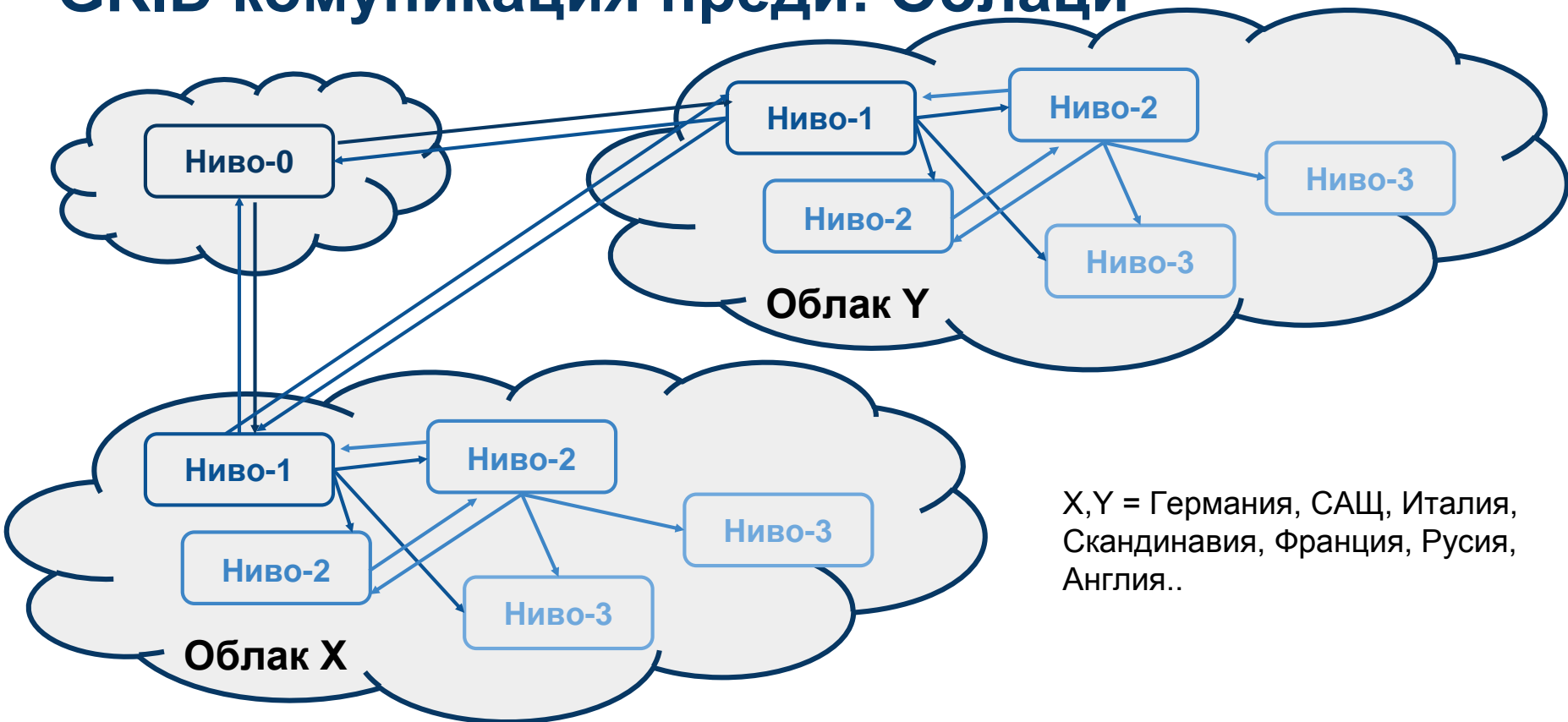
# ATLAS Сайт

- Всеки сайт (обикновено) се състои от ресурси за съхранение на данни + изчислителни мощности
  - Squid proxy cache - за разпространяване на софтуера и “detector conditions” база данни
  
- Нива:
  - Ниво-0 (ЦЕРН)
  - Ниво-1 - Лентови носители + някои услуги (FTS))
  - Ниво-2 - Сайтове с подписан “договор” (Memorandum of Understanding - MoU)
  - Ниво-3 - Непостоянни ресурси

# ATLAS Облак

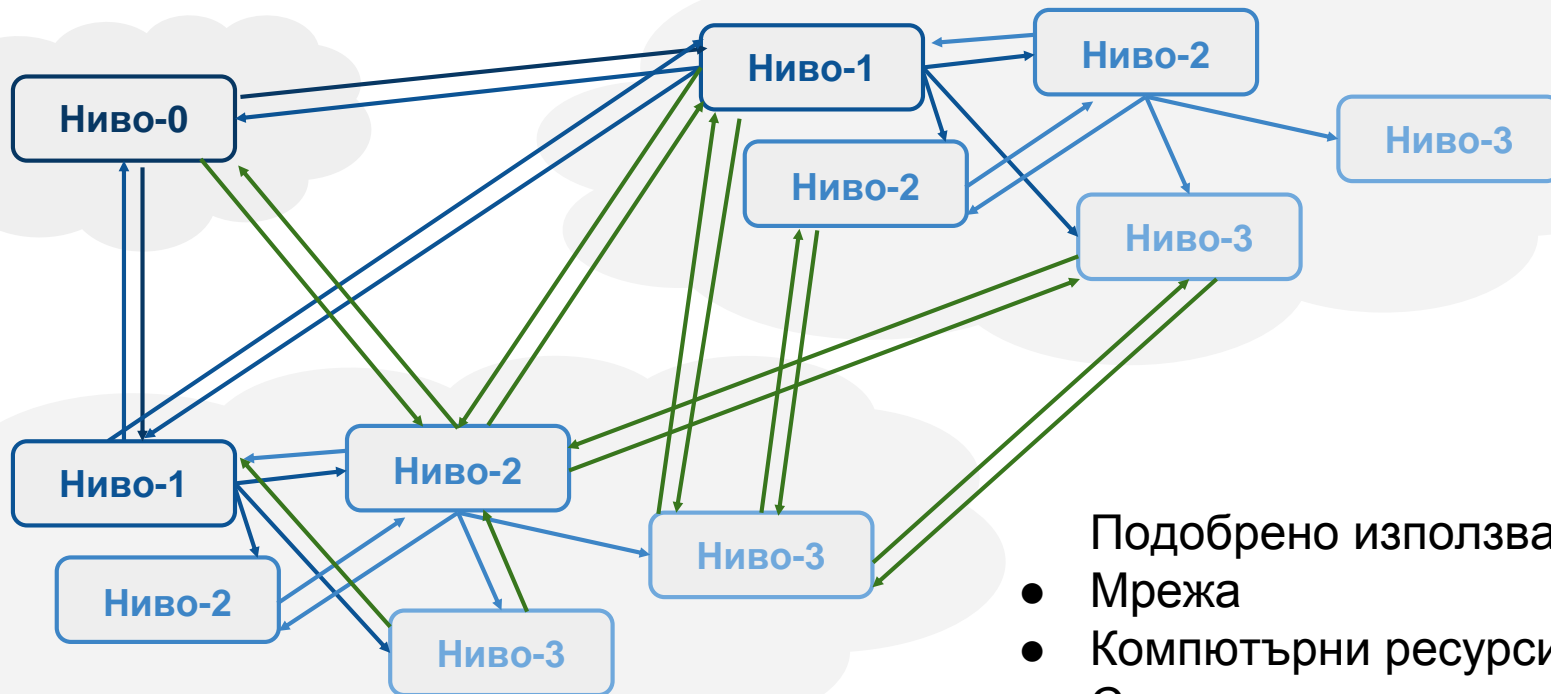
- Няма нищо общо с “облачни изчисления” (AWS, Google Cloud, и т.н.)
- Логическа група от сайтове
  - Един сайт от ниво-1 и няколко ниво-2 и ниво-3
  - Или в един географски регион, или от един източник на финансиране
- Поддръжката е предоставена от екипите на отделните облаци
  - Близо до сайтовете и техните проблеми
  - Най-често - на същия език
- Остаряла концепция
  - Диктувана от ограничения в мрежавите скорости
  - Все още се използва донякъде – за разделяне на екипите за поддръжка.

# GRID комуникация преди: Облаци



X, Y = Германия, САЩ, Италия,  
Скандинавия, Франция, Русия,  
Англия..

# GRID комуникация сега: Тотална мрежа (Full mesh)



Подобрено използване на:

- Мрежа
- Компютърни ресурси
- Съхранение на данни

# Роля на нивата: Ниво-0

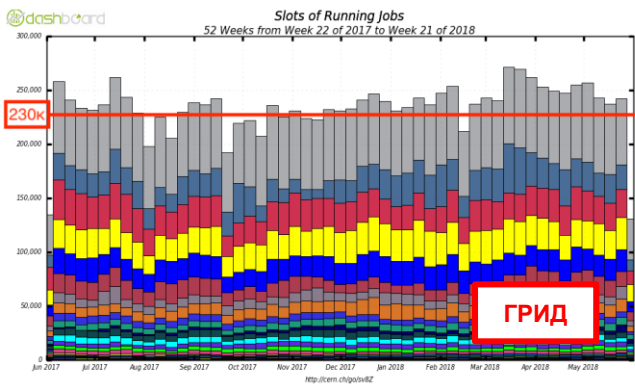
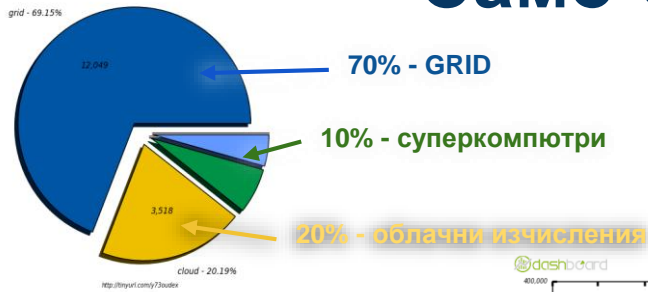
- Ниво-0
  - Копие на всички първични данни на лентови носители
  - Първо ниво на обработка на данните от LHC
  
- Централни услуги в ЦЕРН (стабилност)
  - Разпределяне на задачи - PanDA
  - Разпределяне на данни - Rucio
  - База данни - Oracle
  - Прехвърляне на файлове - FTS (WLCG)
  - Кешова инфраструктура - Frontier
  - Разпределяне на софтуер - CVMFS - Stratum 0 (ниво-0)

# Роля на на нивата: Нива 1, 2, 3

- Нива 1 и 2 - договорени (pledged) компютърни ресурси и дискови и лентови носители
  - Съхранение копие на първичните данни и на вторични данни
  - Разлики между ниво-1 и ниво-2
    - Лентови носители (само в ниво-1)
    - Поддръжка
      - 24x7 за ниво-1
      - В рамките на работния ден - ниво-2
    - Някои ниво-1 сайтове доставят и централни услуги (Frontier, FTS)
      - Близост до сайтовете
      - Сигурност на услугата

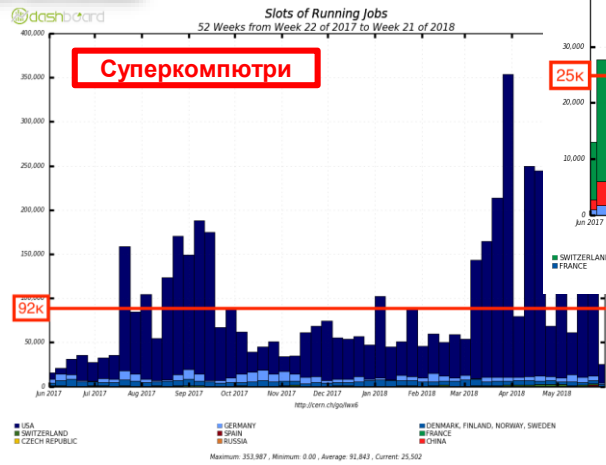
# Симуляции в ATLAS

# Само GRID ли?



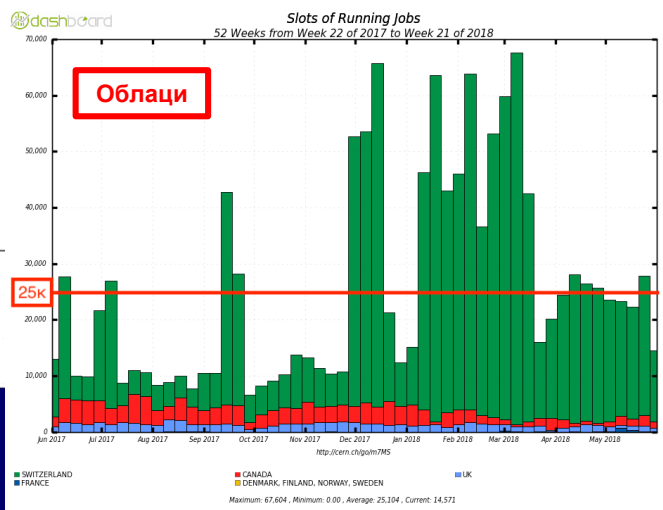
- USA
- GERMANY
- SLOVENIA
- SPAIN
- JAPAN
- ROMANIA
- CHINA
- PORTUGAL
- TURKEY
- SWITZERLAND
- FRANCE
- CANADA
- CZECH REPUBLIC
- TAIWAN
- POLAND
- DENMARK, FINLAND, NORWAY, SWEDEN
- SOUTH AFRICA
- AZERBAIJAN
- UK
- ITALY
- RUSSIA
- NETHERLANDS
- ISRAEL
- AUSTRIA
- SLOVAKIA
- CHILE
- plus 4 more

Maximum: 271.686, Minimum: 0.00, Average: 228.670, Current: 130.673



- USA
- GERMANY
- CZECH REPUBLIC
- SWITZERLAND
- FRANCE
- UK
- ITALY
- RUSSIA
- DENMARK, FINLAND, NORWAY, SWEDEN
- FRANCE
- CHINA

Maximum: 353.987, Minimum: 0.00, Average: 92.842, Current: 25.502



Maximum: 67.604, Minimum: 0.00, Average: 25.104, Current: 14.571

# Разликата

## HTC: High-Throughput Computing (GRID)

- Какво е?
  - Голяма изчислителна мощност за дълго време.
- За какво служи?
  - Ефективното изпълнение на много и слабо свързани задачи.

## HPC: High-Performance Computing (Суперкомпютри)

- Какво е?
  - Голяма изчислителна мощност за ограничено време.
- За какво служи?
  - Ефективно изпълнение на много, тясно свързани задачи.

## Облачни изчисления

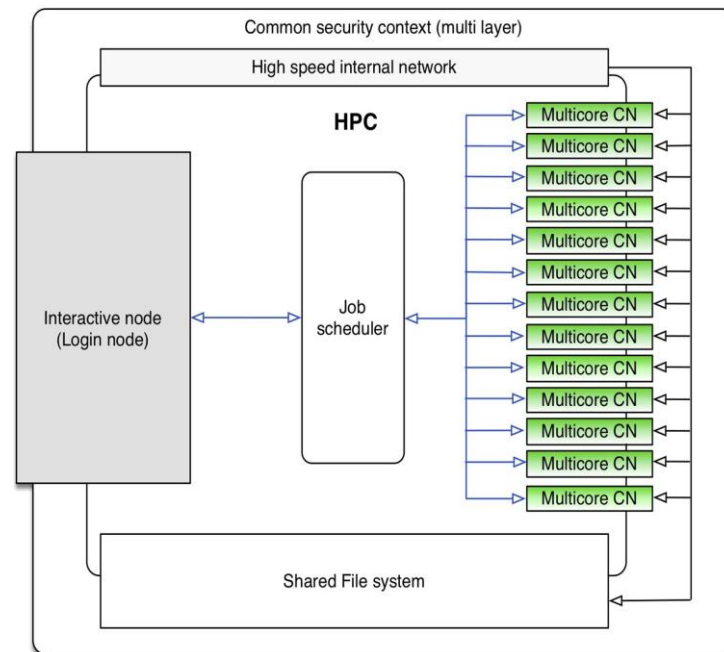
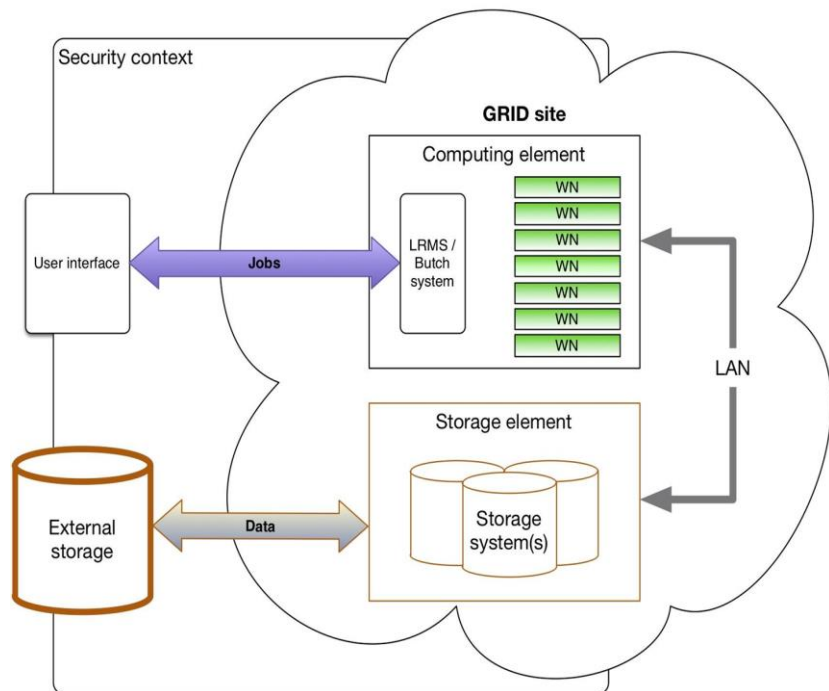
- Какво е?
  - Виртуални машини - колкото и каквито са нужни
- За какво служи?
  - Рарантирани ресурси в сигурна и изолирана среда без нужда от поддръжка на хардуер. API достъп



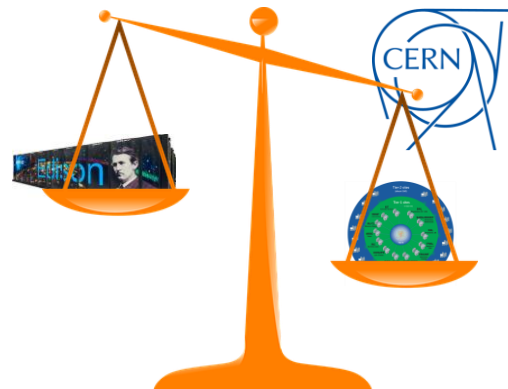


# Суперкомпютри

# GRID и Суперкомпьютри



# GRID и Суперкомпютри



- GRID се състои от компютърни клъстери
- По-голяма част от изчисленията в другите области на науката стават в суперкомпютри.
- Разлики: 

	суперкомпютри <>	GRID:
○ Нишки:	много (~100 000 нишки)	една
○ Работа на отделния компютър:	част от задача	една задача
○ Входно - изходни операции:	малко	много
○ Брой на файлове:	малко	много
○ Потребителска идентификация:	логин/парола (еднократна!)	сертификат
○ Процесори / операционни системи:	много	няколко
○ Директна връзка с интернет:	не	да
○ Местонахождение	едно	много

# Суперкомпютри използвани от АТЛАС

## Най-големите...



[Titan](#) (САЩ, [top500: 7](#))  
Cray XK7  
Opteron 6274 16C 2.2 GHz  
560640 ядра / 694 TB памет / 40 PB диск  
Производителност: 27.11 PFLOP/сек

## И още много в...

САЩ, Германия, Швейцария, Чехия,  
Испания, Франция, Китай, Скандинавия,  
Русия..



[Theta](#) (САЩ, [top500: 21](#))  
Cray XC40  
Intel Phi 7230 64C 1.3GHz  
280320 ядра / 914 TB памет / 11 PB диск  
Производителност: 11.66 PFLOP/сек

## Еднаквото между всички..

Няма такава..



[Edison](#) (САЩ, [top500: 104](#))  
Cray XC30  
Intel Xeon E5-2695v2 121C 2.4GHz  
133824 ядра / 357 TB памет / 7.56 PB диск  
Производителност: 2.57 PFLOP/сек

.....

# Мащаб на задачите на един суперкомпютър

Вид задача	Минимален брой компютри	Максимален брой Компютри	Максимална дължина на задачата	Приоритет
1	11,250	—	24.0	15
2	3,750	11,249	24.0	5
3	313	3,749	12.0	0
4	126	312	6.0	0
5	1	125	2.0	0

# HPC: Backfill



# HPC: Backfill



Източник: David Cameron

# Суперкомпютри: Backfill

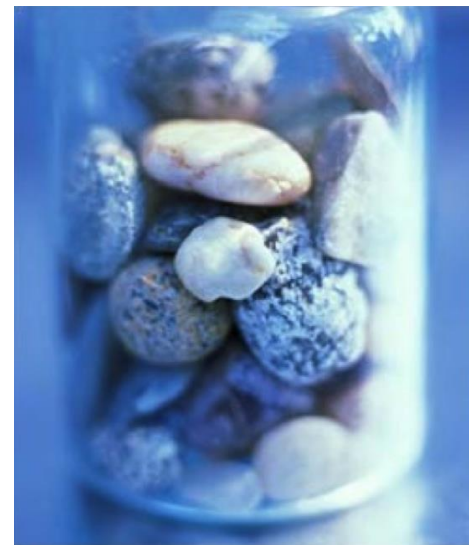


Източник: David Cameron



# Суперкомпютри: Backfill

(пълним дупки всекакви)



Източник: David Cameron

# Суперкомпютри: Backfill

(пълним дупки всекакви)



Източник: David Cameron

# Суперкомпютри: Backfill

(пълним дупки всекакви)

## Размер на дупките

- 400M CPU\*часа на година (Titan)
- Обичайна ефективност на използване на суперкомпютър във времето - 90%.
- 10% == 400M CPU\*часа на година (Titan)
- Всяка неефективност е нежелана!

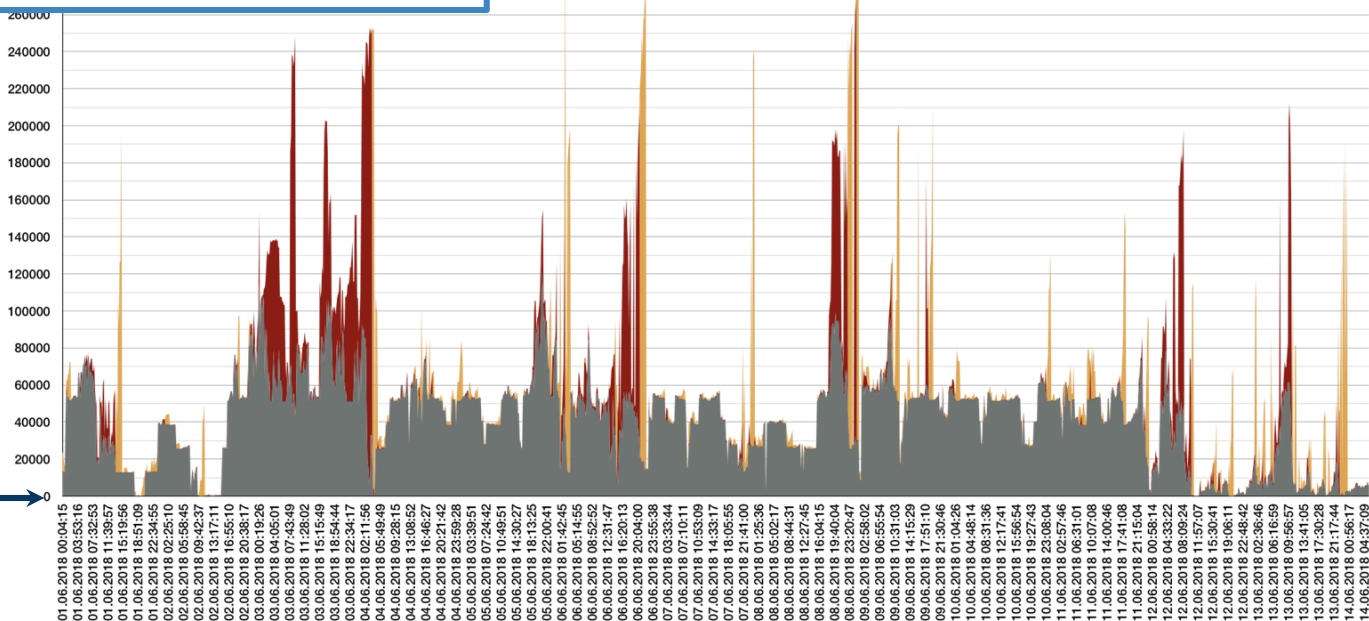
## Запълване

- С какво?
  - Много, с различен размер, но малки задачи (т.е. - ние)
- На каква цена?
  - Безплатно
- И на края..
  - Всички живели щастливо..



# Използване на Titan от ATLAS

- Изчисления по квота
- **Backfill**
- **Свободни ресурси**

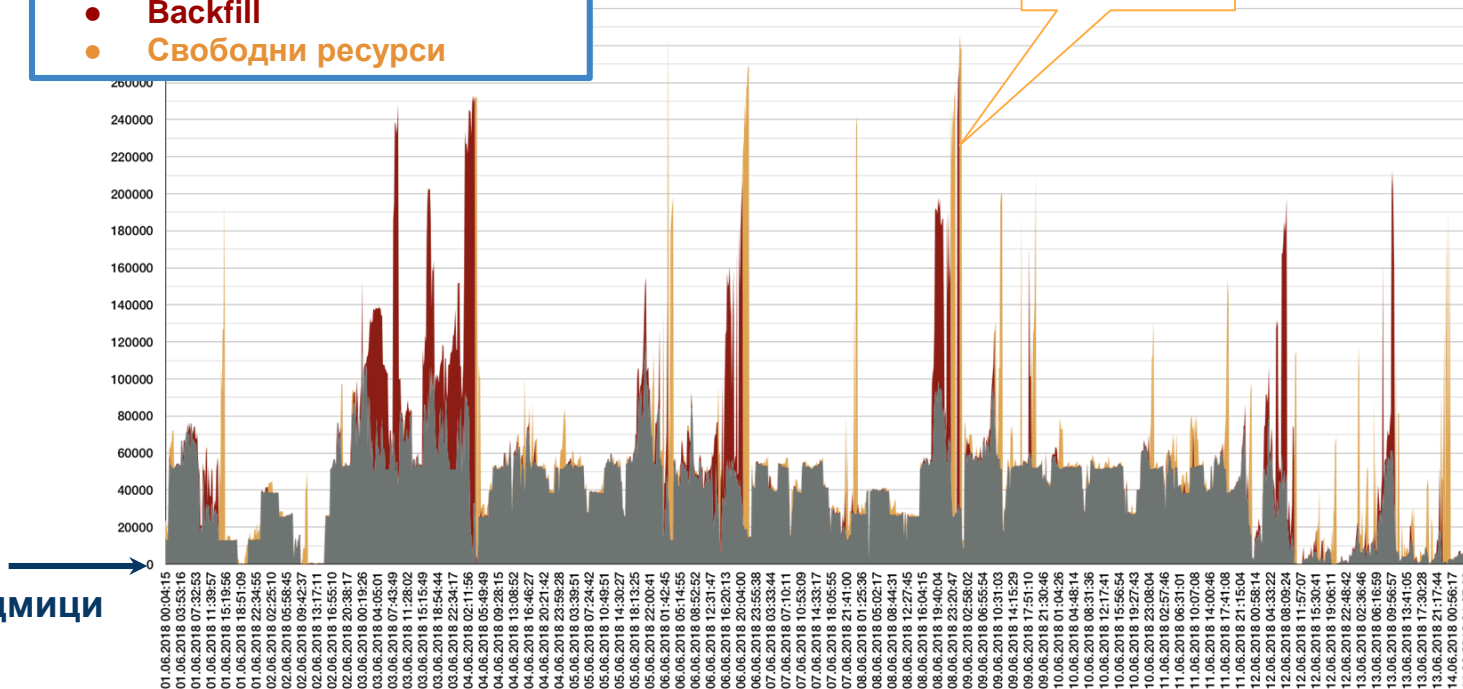


две седмици

# Използване на Titan от ATLAS

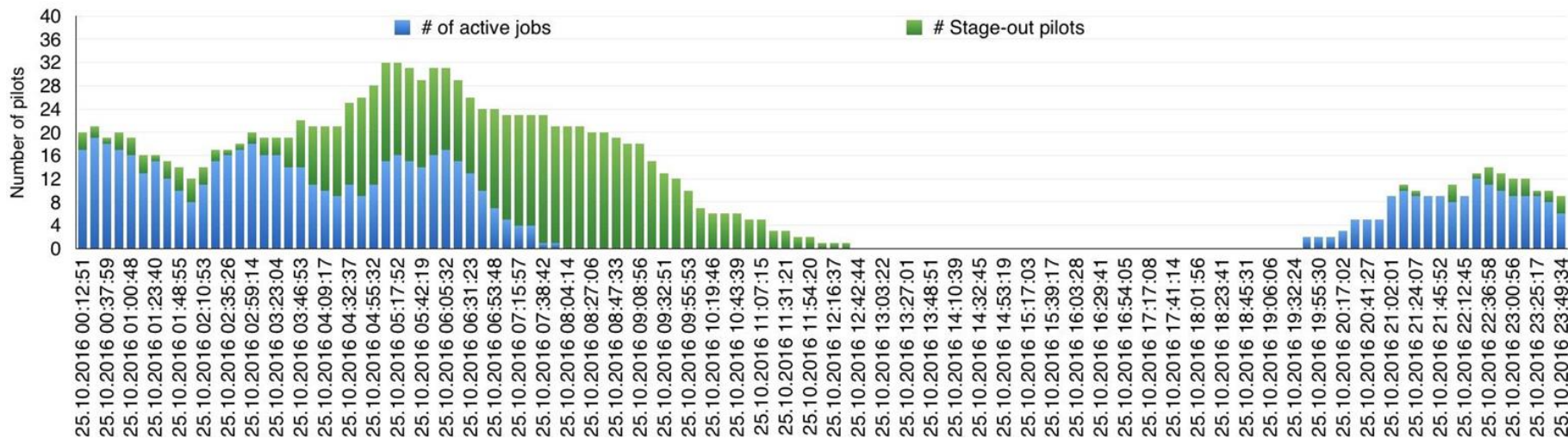
- Изчисления по квота
- **Backfill**
- **Свободни ресурси**

ЗАЩО?!



две седмици

# Проблемът





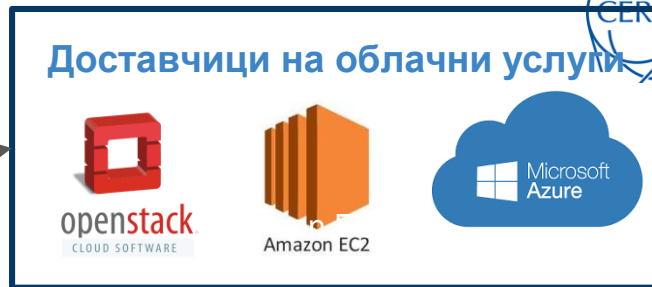
# Облаци

# На кратко

- ATLAS използва собствени, научни и комерсиални.
- Предимства
  - Комерсиални: При нужда от много ресурси за кратко време
  - Научни: Споделяне на ресурси и цена между различни институти и различни области на изследване
  - Собствени: Бързо и контролирано усвояване на ресурси
- Недостатъци
  - Комерсиални: Скъпи



# Принцип на работа

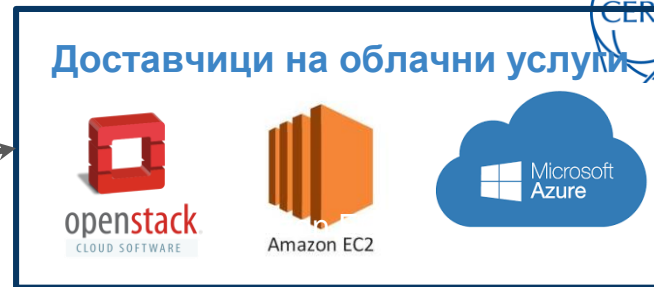


Изпраща задачи

# Принцип на работа



Нарави ми виртуални машини



Наблюдава опашката от задачи

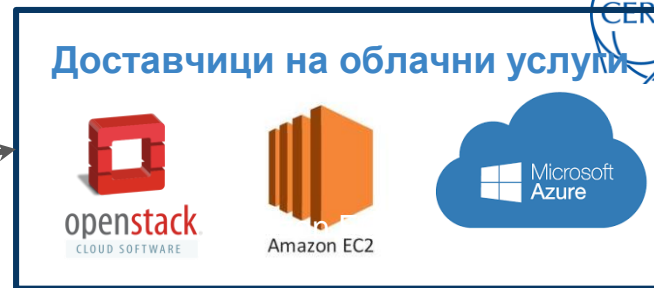


Изпраща задачи

# Принцип на работа



Нарави ми виртуални машини



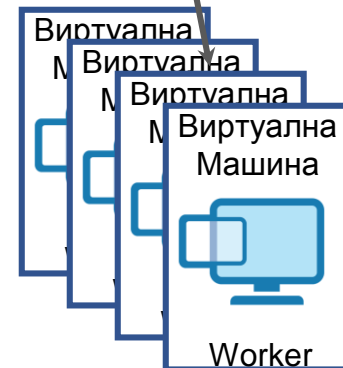
Наблюдава опашката от задачи



създаване



Изпраца задачи

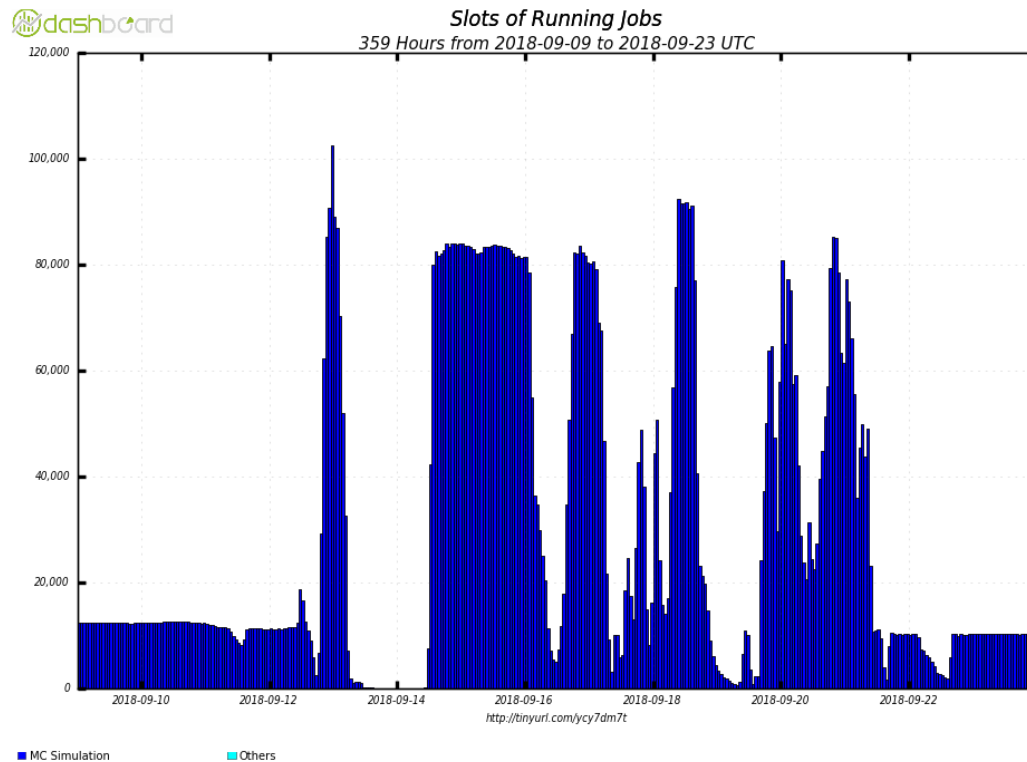


# Принцип на работа



# Собствени облаци - тригерна ферма

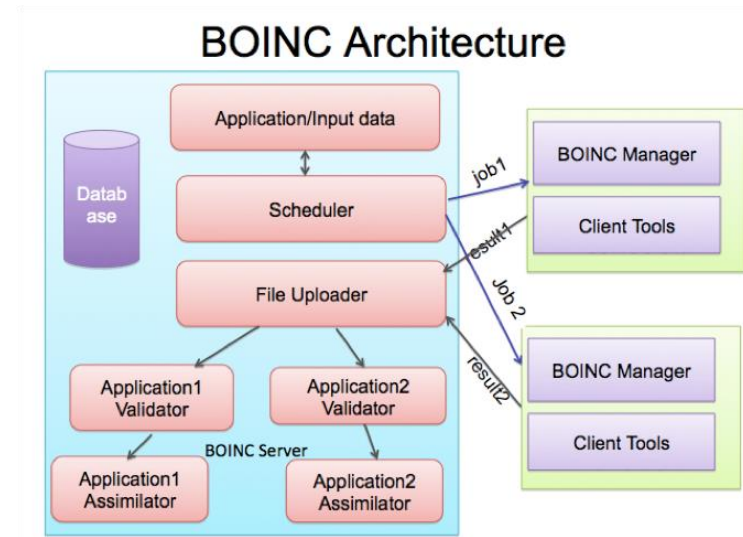
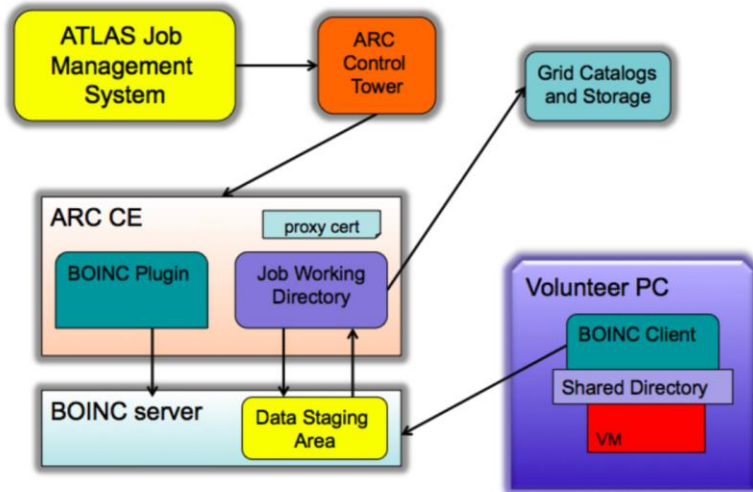
- Ресурс - 120 000 ядра
- Използване - само когато не се набират данни
- Облачно решение:
  - Бързо усвояване на ресурси
  - Няма нужда от инсталиране
  - Светкавично връщане на ресурси при нужда (kill)
- Все още не се използва максимално



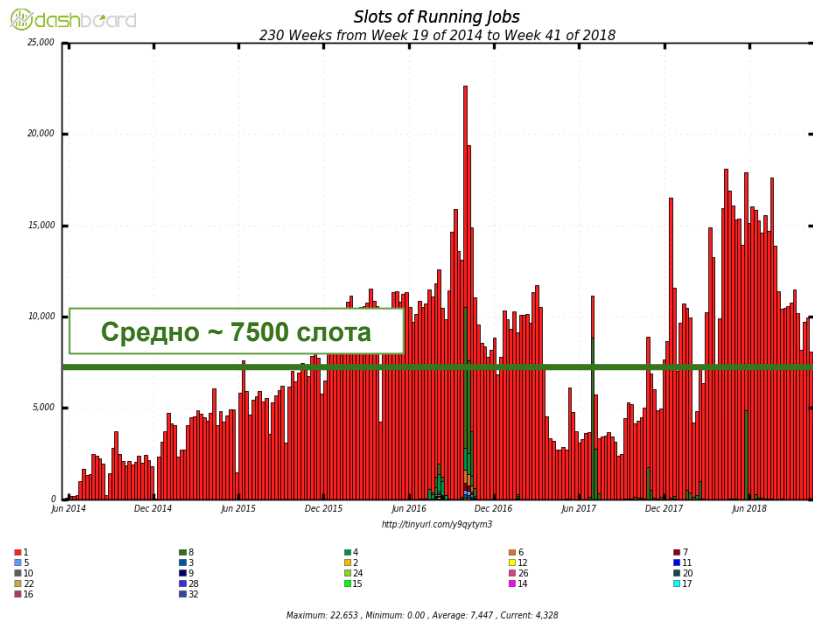


# **Доброволчески споделени изчисления (Volunteer Computing)**

# BOINC



# ATLAS@Home

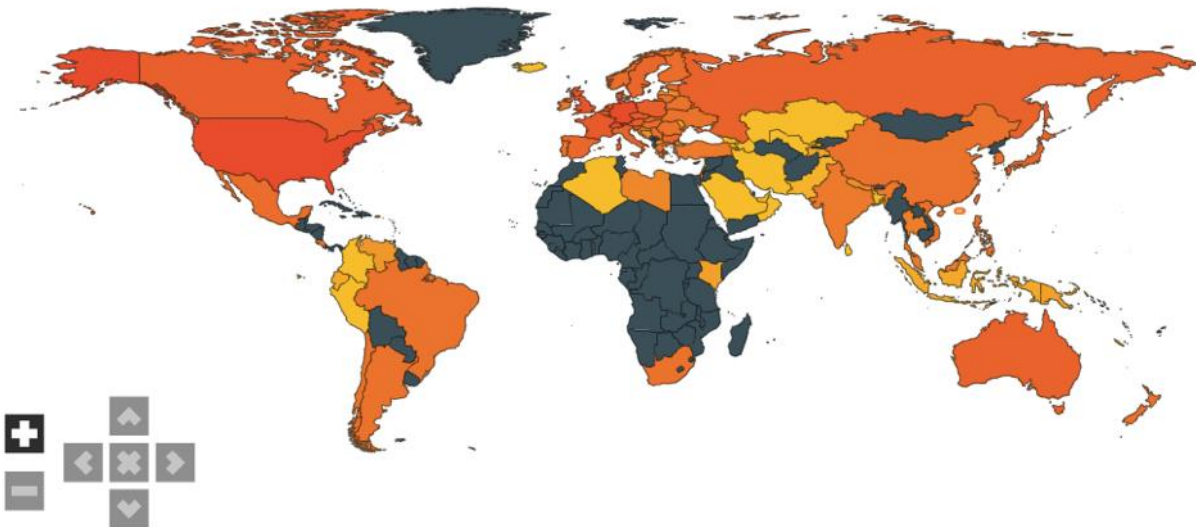
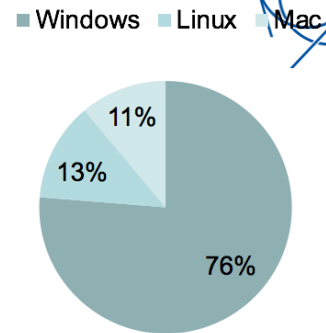


- Обработка ~ 1% от симулираните събития в ATLAS
- Сравним с голям безплатен компютърен център
- Не точно.. доброволците очакват все пак някаква поддръжка
- Голям потенциал в настолните компютри по учреждения които не се изключват нощно време



# Кой ни помага

- VOINC е огромно общество от отдадени хора
- Много от тях участват в няколко проекта едновременно



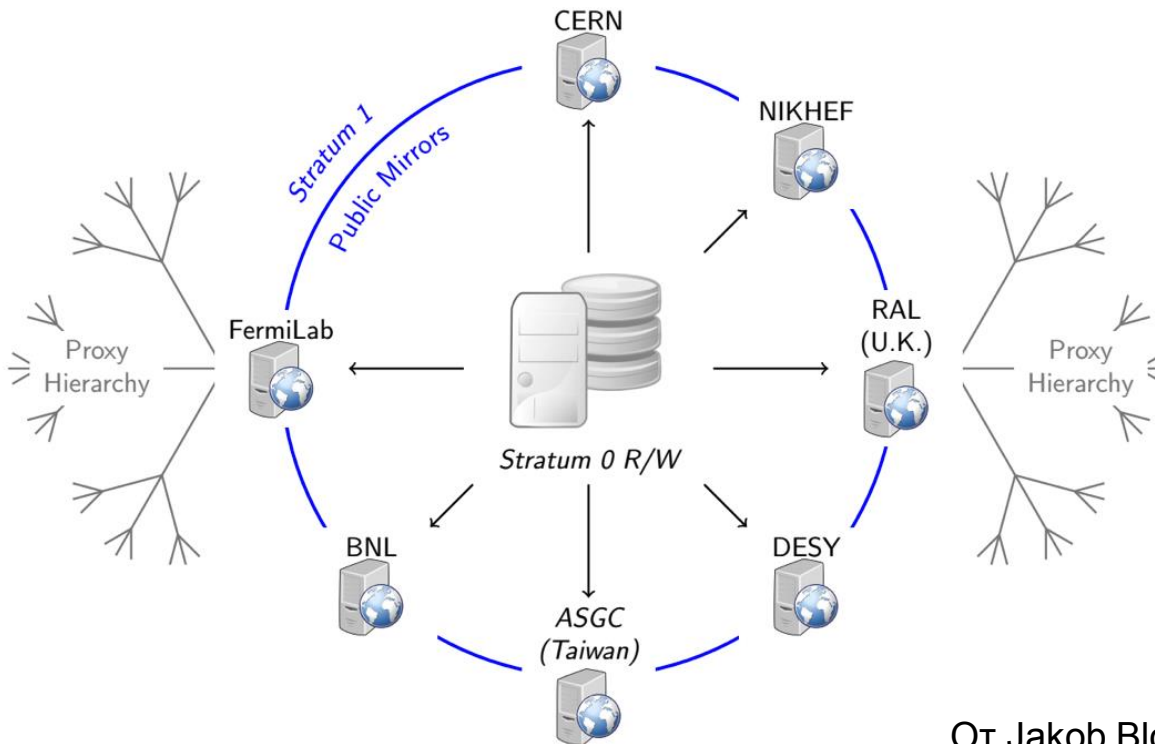
# И ние инсталираме VOINC

- На новопристигналите машини в ЦЕРН
- На машини с “неефективни” услуги
- На малки GRID сайтове
- На сайтове с ограничения



# Дистрибуция на софтуера: CVMFS

- CVMFS - CERN VM FileSystem
  - Мрежова файлова система базирана на http и оптимизирана да доставя софтуера на експериментите по сайтовете бързо и сигурно
  - Новия софтуер се поставя на Stratum-0 и автоматично се репликира на всички Stratum-1
  - Най-голямата част от репликираното се поема от кешовете (Squids) на всеки един сайт
    - Компютрите на всеки сайт четат софтуер само от локалния кеш
    - В случай че локалния кеш не работи, компютрите четат от следващото ниво
- Всички стандартни сайтове в АТЛАС използват CVMFS
  - Трябва връзка с външния свят - не работи при повечето суперкомпютри



От Jakob Blomer

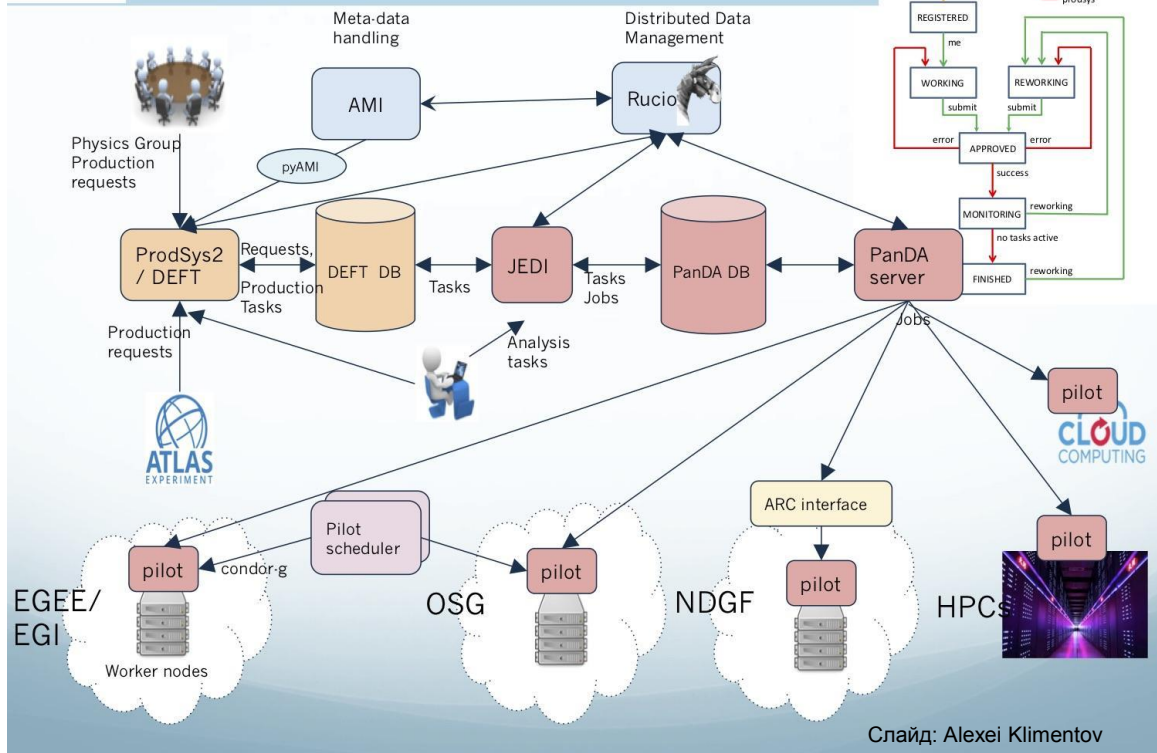


# Изчисления (как)

# Изчислителни задачи: от А до Я



## ATLAS Workflow and Workload Management



## Дефиниции

### Разпределени изчисления:

- Задача (Job) - какво трябва да се сметне на един "компютър"
- Пилот - малка програма която се изпраща на сайта и вика реалната задача от системата за задачи
- HPC - Суперкомпютър

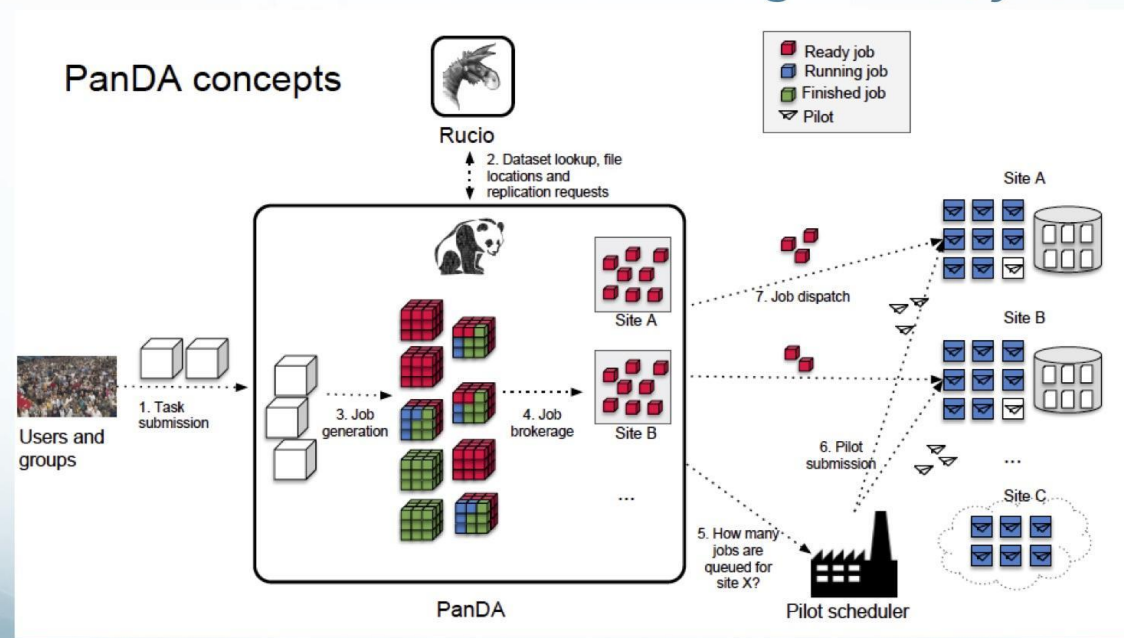
### АТЛАС:

- ProdSys2 / DEFT - системата която взема една заявка (request) и я разпределя
- Rucio - системата за данни - записване, изтриване, преместване, репликиране и т.н.
- Задание (Task) - какво трябва да се пресметне от данни получени пре еднакви условия

# Система за разпределяне на задачи в АТЛАС



## PanDA Workload Management System



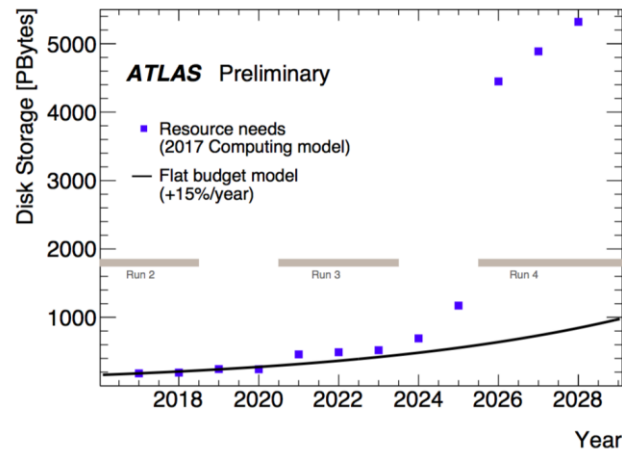
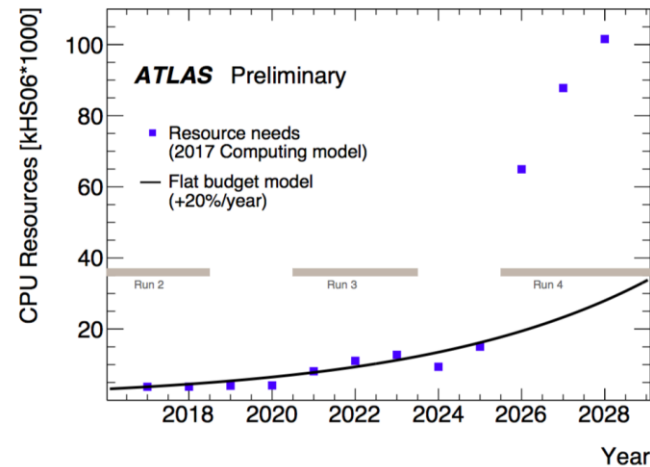
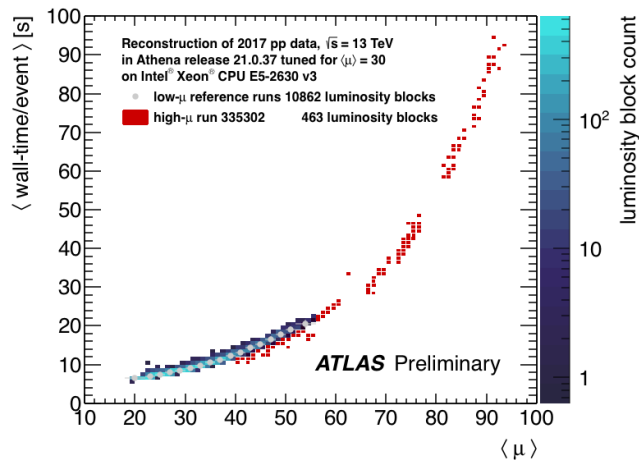
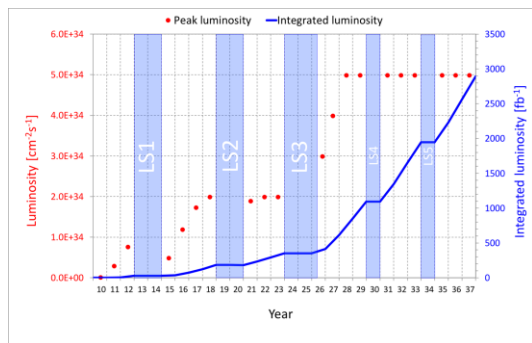
CHEP talk about new PanDA developments: "Harvester: an edge service harvesting heterogeneous resources for ATLAS",  
 Track 3 – Distributed computing, 12 Jul 2018, 11:15  
<https://indico.cern.ch/event/587955/contributions/2937391/>  
 Alexei Klimentov





# Бъдещето

# HL-LHC



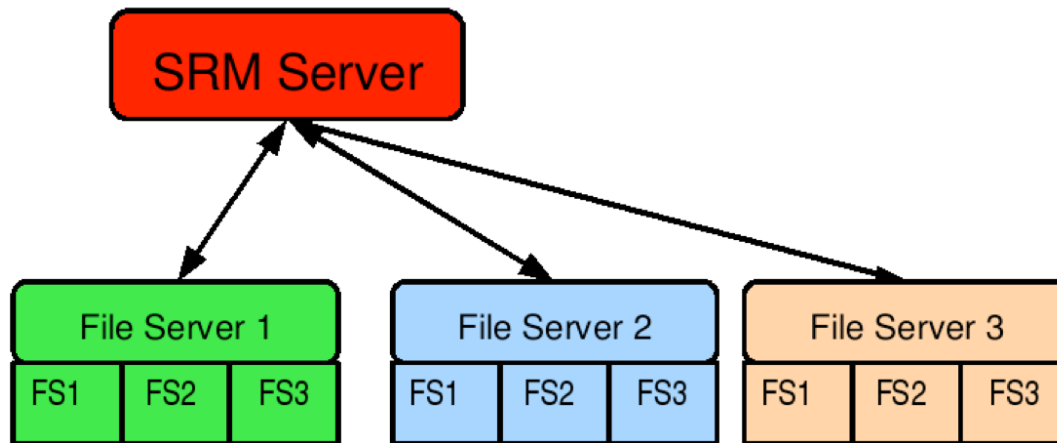


**Въпроси?**  
**Ivan.Glushkov@cern.ch**

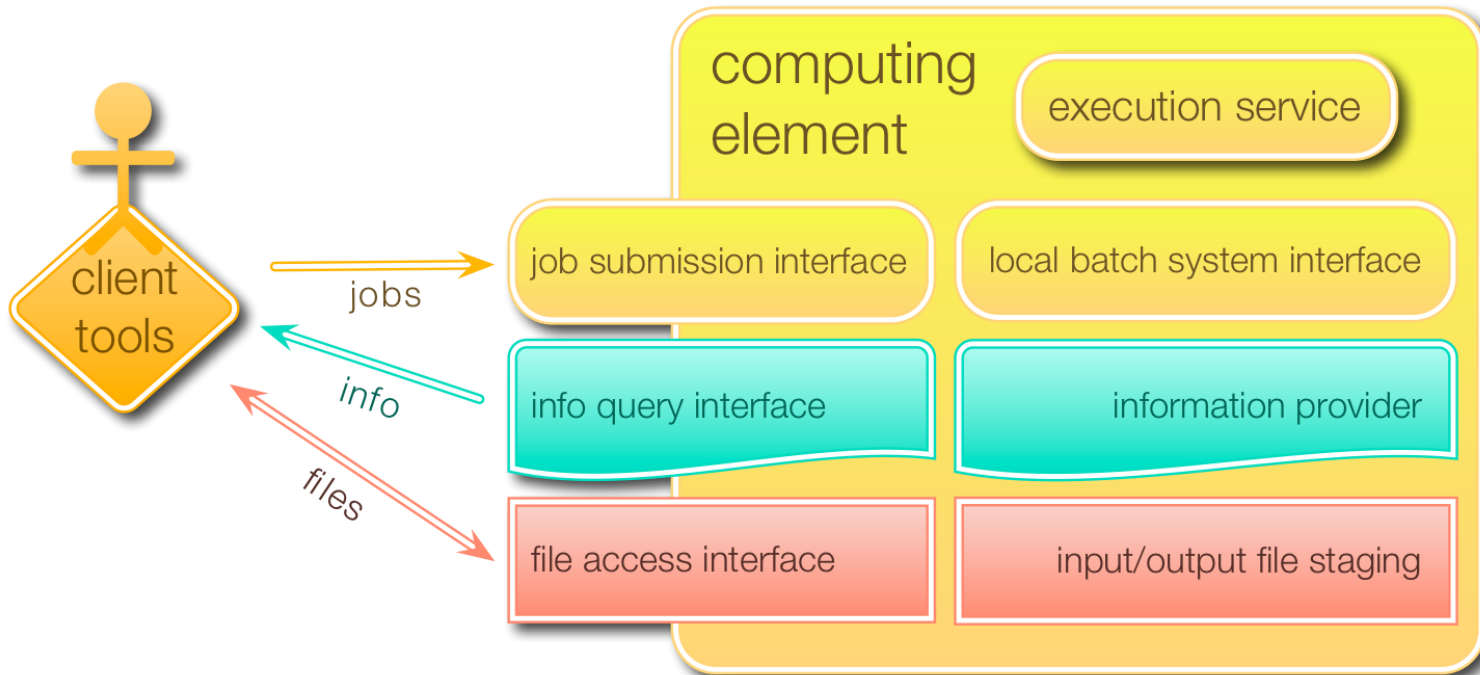


# Подробности

# SE (Storage Element)



# CE (Computing Element)



# Global Shares

