

Тенденции в информационной инфраструктуре ядерно-физических центров

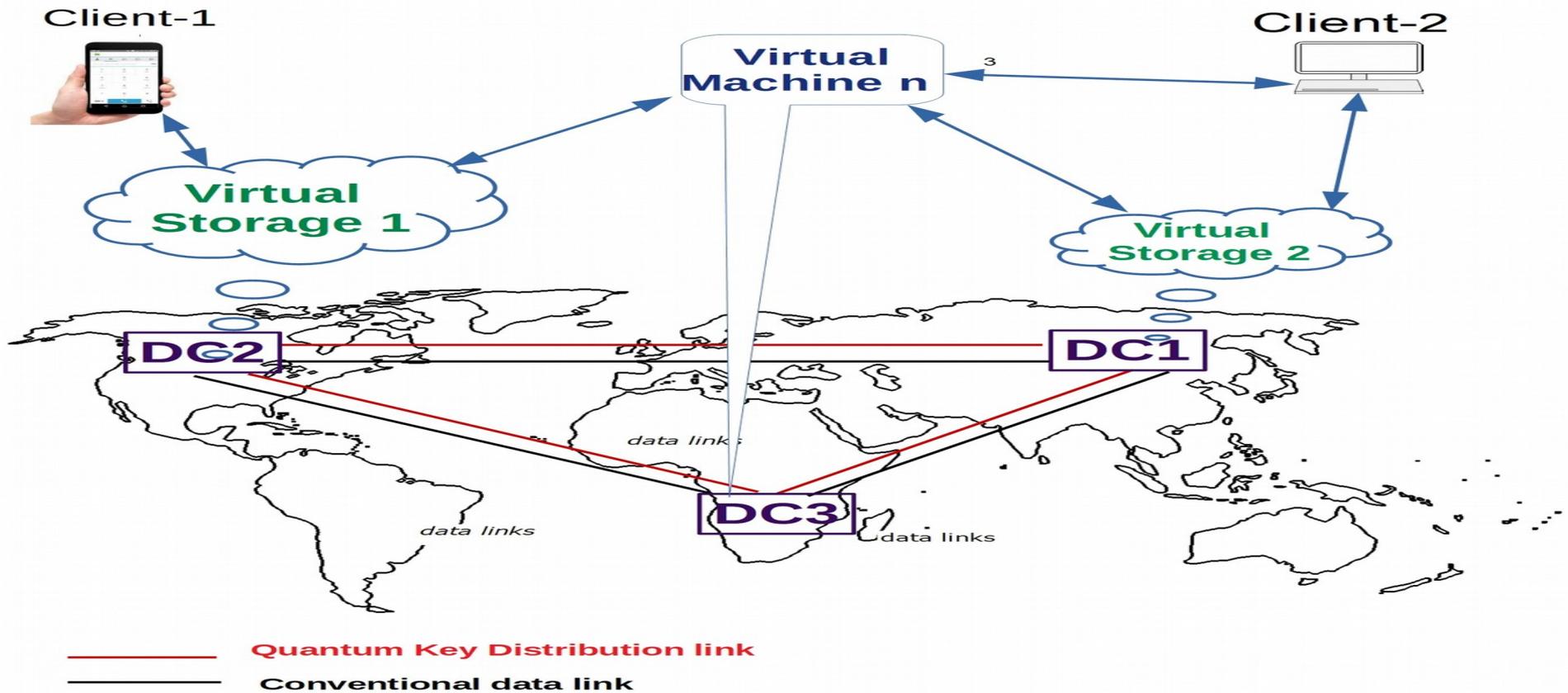
Семинар ОФВЭ

Докладчик: Андрей Е Шевель

План семинара

- Введение
- Информационная инфраструктура
- Параметры информационной инфраструктуры в исследовательских центрах
- Взаимодействие с информационной инфраструктурой
- Доступность, сохранность, безопасность
- Технологические изменения

Завершение проекта No.03.G25.31.0229



Информационная инфраструктура

- Совокупность организационных практик, технической инфраструктуры, социальных коммуникаций, которые в совокупности реализуют эффективный обмен данными между исследователями, что ускоряет получение результатов.

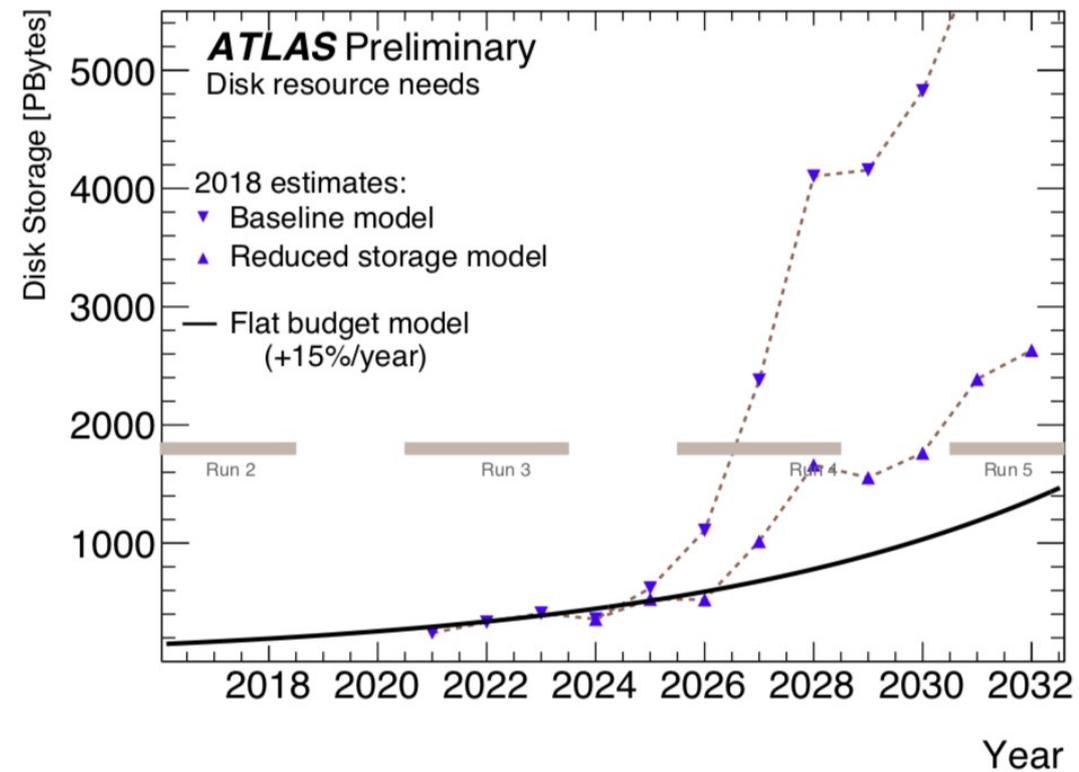
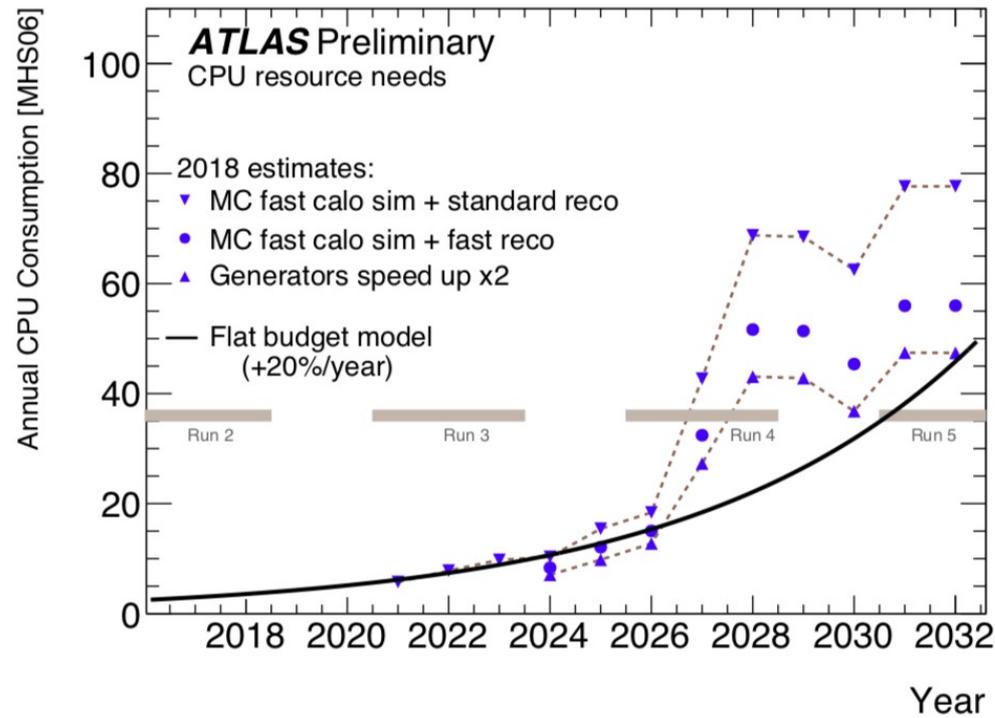
Большие переменны

- **Большие данные**
 - Big data: **Velocity, Variety, Volume**
 - Большие вычислительные мощности
 - Ёмкость каналов передачи (100 Gbit и более)
- **[Компьютерные] облачные системы**
 - On-demand self-service, Broad network access, Resource pooling, Rapid elasticity, Measured service
 - **Облачный стиль**
 - Компьютерная инфраструктура
 - Астрономические наблюдения
 - Другие применения
- **Нейронные сети**

Данные

- Данные или информация ?
- Основные операции с данными
 - Хранение данных
 - Передача данных
 - Преобразование данных
 - Визуализация данных

Пример из АТЛАС



CERN

- Число процессорных ядер $> 240\text{K}$;
- Число дисководов $> 70\text{K}$ (около 500 PB);
- Число магнитофонов в роботизированном хранилище ~ 76 ;
- Число сотрудников в ИТ подразделении ~ 430 .

IT personnel ~ 432 [<https://cds.cern.ch/record/2677223/files/CERN-HR-STAFF-STAT-2018-RESTR.pdf>]

IN2P3

- **Resources**

- 65 IT engineers
- budget : 7M€

- **Computing**

- 1000 servers, 40k threads, 440kHS06

- **Storage**

- Tapes : 70 PB
- HDDs : 30 PB

- **Networks**

- 20 Gbps to Renater
- 40 Gbps to CERN
- 40 Gbps to LHCONE/WLCG

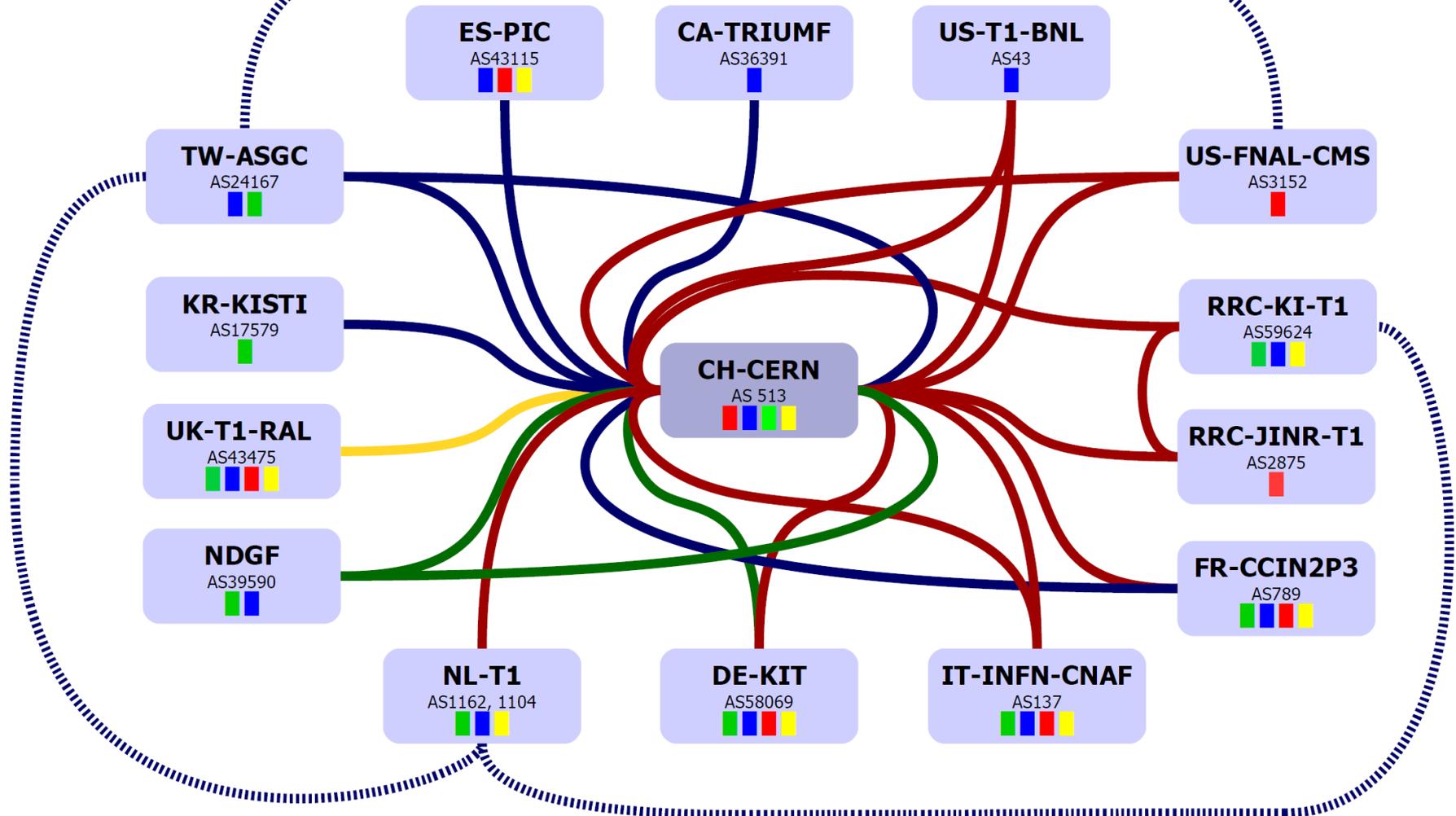
Параметры SDCC (BNL)

- ▶ RACF occupies ~15,000 sq. ft.
- ▶ ~250 Computer Racks + 9 Robotic Linear Magnetic Tape Libraries (~85k tape cartridge cells, ~200 PB Capacity)
- ▶ HPSS based tape management system, one of the largest instances worldwide
- ▶ Scalable parallel filesystem with GPFS (currently deployed instance scales to 100 GB/s)
- ▶ 50,000 X86 based compute cores, 125 TB memory, 35 PB Disk
- ▶ Bisection bandwidth of 3 Terabits/sec between storage and compute resources
- ▶ ~800 kW Power Utilization (~7 MWh per year)

Магистральные компьютерные сети

- LHCOPN
 - LHC Optical Private Network
- LHCONE
 - LHC Open Network Environment
- ESnet
 - Energy Science Network

LHCOPN

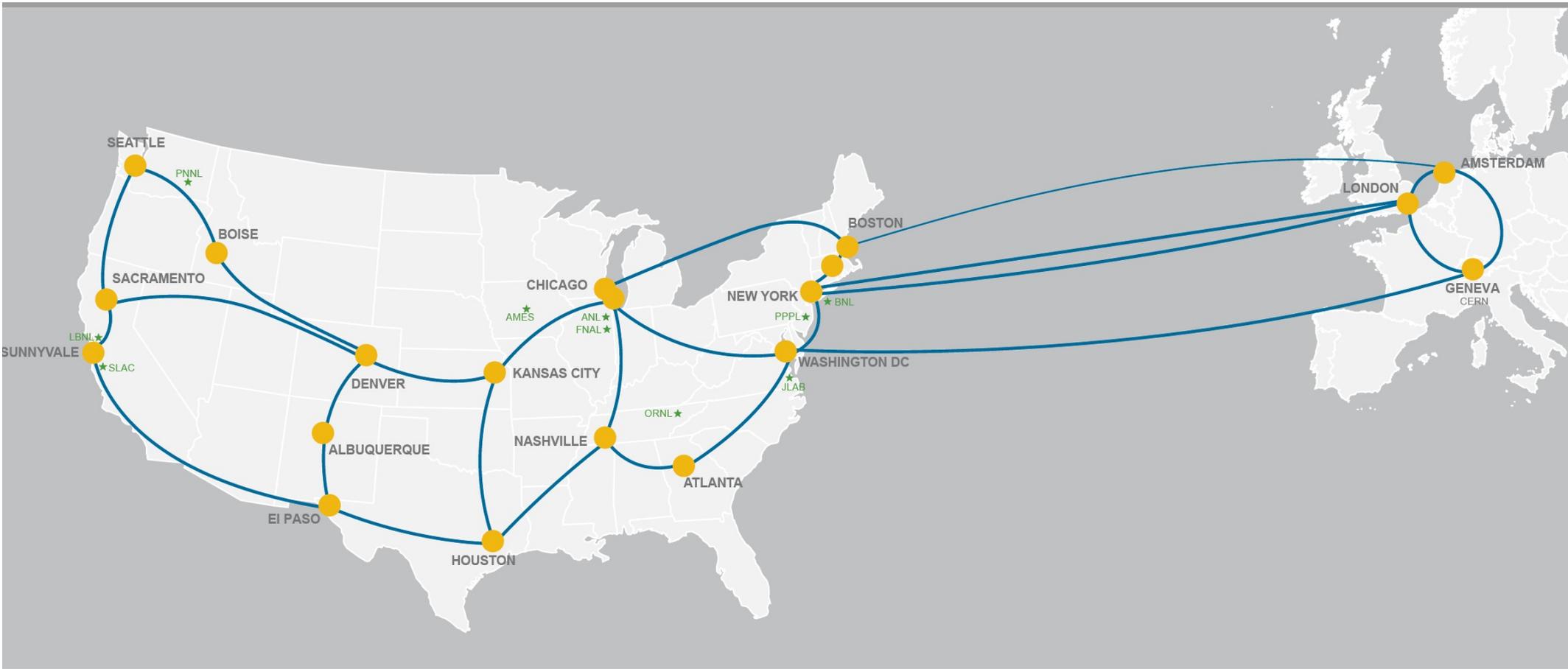


Ref: <https://twiki.cern.ch/twiki/bin/view/LHCOPN/OverallNetworkMaps>

	T0-T1 and T1-T1 traffic		10Gbps
	T1-T1 traffic only		20Gbps
	= Alice		30Gbps
	= Atlas		40Gbps
	= CMS		100Gbps
	= LHCb		

edoardo.martelli@cern.ch 20190516





ESnet

ENERGY SCIENCES NETWORK

★ Department of Energy Office of Science National Labs

- Ames** Ames Laboratory (Ames, IA)
- ANL** Argonne National Laboratory (Argonne, IL)
- BNL** Brookhaven National Laboratory (Upton, NY)
- FNAL** Fermi National Accelerator Laboratory (Batavia, IL)
- JLAB** Thomas Jefferson National Accelerator Facility (Newport News, VA)

- LBNL** Lawrence Berkeley National Laboratory (Berkeley, CA)
- ORNL** Oak Ridge National Laboratory (Oak Ridge, TN)
- PNNL** Pacific Northwest National Laboratory (Richland, WA)
- PPPL** Princeton Plasma Physics Laboratory (Princeton, NJ)
- SLAC** SLAC National Accelerator Laboratory (Menlo Park, CA)

Нерядкие особенности в исследовательских центрах

- Все виды компьютерной инфраструктуры входят в сферу ответственности ИТ департамента, включая
 - Телефонную связь;
 - Все виды десктопов и операционных систем;
 - Все виды компьютерных сервисов (обновления, безопасность, сетевая инфраструктура, резервное копирование, проч).
 - ИТ департамент имеет свой бюджет.

Jupyter

Applications Places Google Chrome en Thu 10:37

What's New at the RACF x SDCC JupyterHub Home x +

← → ↻ ↵ jupyter.sdcc.bnl.gov ☆ ⋮ 👤 ⋮

Apps Linux Container... Новая газета -... @ :: General Purp... OpenStack Sum... Индекс Хирша... ПК HP - Клави... Как подключит... Как подключит... BEFA-Science New Site Comp... » | Other bookmarks

SDCC JupyterHub

The SDCC offers multiple JupyterHub instance and back-end combinations for different users and accounts. Choose the appropriate option from the instances displayed below.

[More information](#)

[Questions and support](#)



SDCC **HTC**

Access to Condor queues and HTC computing resources via SDCC JupyterHub. Requires a valid SDCC account and corresponding experiment affiliation.

[Launch](#) [More info](#)

SDCC HTC JH



SDCC **HPC**

Access to Slurm scheduling and GPU computing resources on the IC and KNL clusters via JupyterHub. Requires a valid SDCC account and computing resource allocation.

[Launch IC](#) [Launch KNL](#) [More info](#)

SDCC HPC JH

REANA

- REANA (Reusable Analyses) - <https://cds.cern.ch/record/2652340/files/Fulltext.pdf>

Платформа REANA состоит из набора микросервисов, позволяющих запускать и отслеживать в облаке рабочие задания вычислительных процессов с использованием контейнеров.

- Позволяет использовать результаты рабочего процесса, запустив интерактивный сеанс в ноутбуке Jupyter.
- Можно использовать из Gitlab CI/CD (continuous integration and delivery).

Доступность сохранённых данных

The New York Times

Archive

WORLD U.S. N.Y. / REGION BUSINESS TECHNOLOGY SCIENCE HEALTH SPORTS OPINION ARTS STYLE T

New York Times Article Archive

1851–PRESENT

The complete archive of The New York Times can now be searched from NYTimes.com — more than 13 million articles total.

Searching the Archives

The archive is divided into two search sets: 1851–1980 and 1981–present.

[Search the Article Archive: 1981-Present »](#)

[Search the Article Archive: 1851-1980 »](#)

 FACEBOOK

 TWITTER

 GOOGLE+

 SAVE

 EMAIL

 SHARE

 PRINT

 REPRINTS

Оцифровка мультимедиа в CERN

- Уже оцифровано:
 - 2500 видео лент;
 - 1000+ персон опознано на видео и негативах;
 - несколько тысяч фото негативов,
 - много сотен аудио лент
- И загружено на сервер CDS.cern.ch;

Сохранность данных во времени

- Сколько времени хранить: 20? 30? 50 лет или больше ?
- Перемещение данных с устройств старого образца на устройства хранения новых типов
 - Когда речь идёт о сотнях РВ или единицах ЕВ – процесс перезаписи становится постоянной работой с постоянным штатом сотрудников.
- Программы обработки данных (или симуляции) необходимо поддерживать чтобы можно было проверить и/или уточнить алгоритмы много десятков лет спустя ...

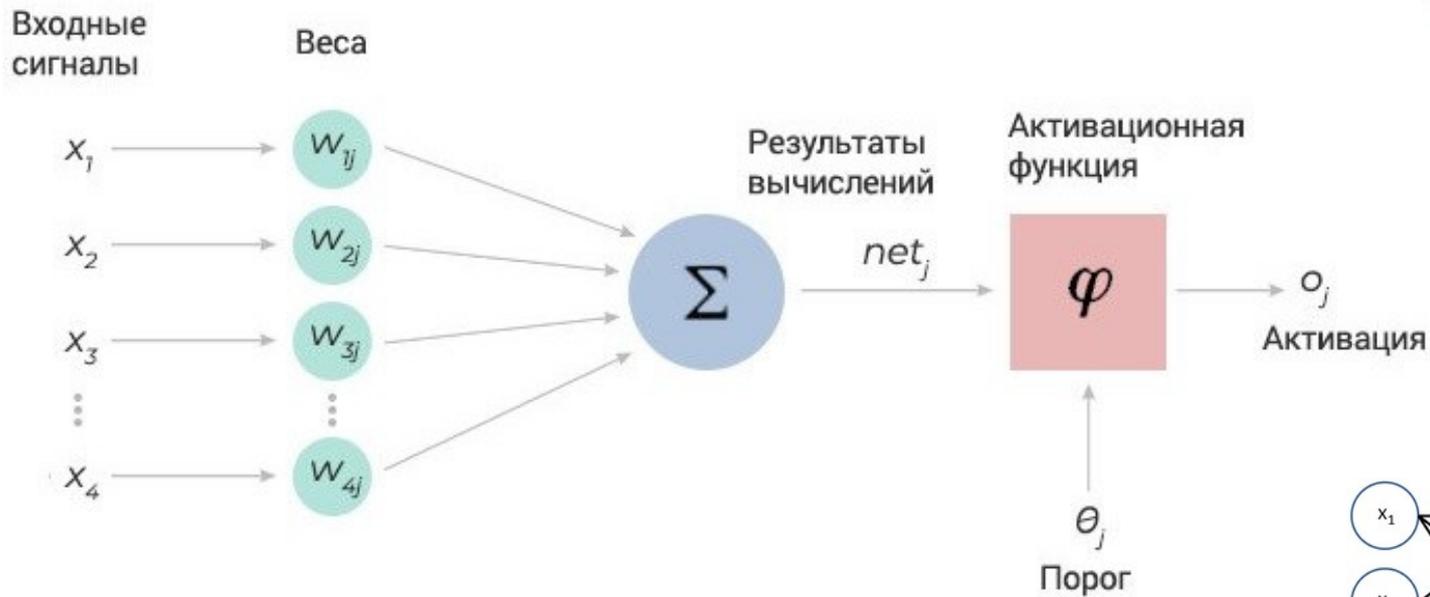
Воспроизводимость

- Научные результаты должны быть воспроизводимыми.
- Результаты измерений и обработки должны быть открытыми, свободно доступными в том числе для поиска.
- Сложные рабочие процессы, записанные в электронных ноутбуках, должны адаптироваться к новым потокам данных и новому программному и аппаратному обеспечению.

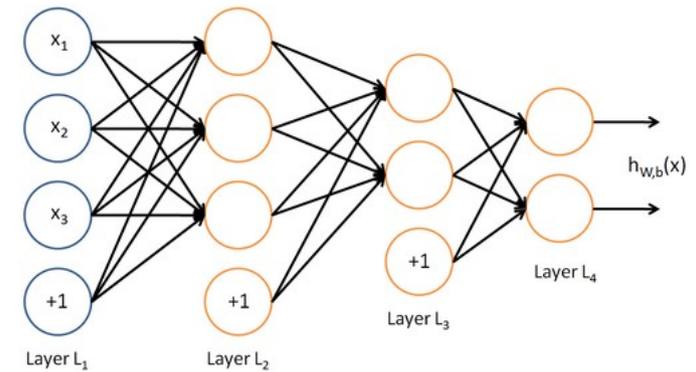
Новые угрозы

- Безопасность (несанкционированный доступ к данным и ресурсам)
 - Персональные данные
 - Большой вопрос “что есть персональные данные?”
 - Московский суд установил, что ваши фото и видео отдельно от других сведений о вас не являются персональными данными.
 - FZ-152 (RU), GDPR - The GDPR aims primarily to give control to individuals over their personal data and to simplify the regulatory environment for international business by unifying the regulation within the EU.
 - Свои данные на своих серверах?
 - Мейл сервер ЦЕРН
 - Новые правила создают и новые проблемы (новый термин: security cratia).
- Мы живём в прозрачном мире.
- Опасность объединения хранилищ данных.

Нейронные сети



Простейшая нейронная сеть



И чуть сложнее

Frank Rosenblatt создал персептрон (первая нейронная сеть) в 1959.

Технология WiFi

- "Wi-Fi mostly" - пилотная инициатива, пока для ИТ подразделения CERN:
 - Мобильным устройствам не рекомендовано подключаться к структурированной проводной сети;
 - В результате более сотни сетевых розеток вообще не используются;
 - Пока в CERN не отработана загрузка Линукс с использованием WiFi.
 - https://indico.cern.ch/event/810635/contributions/3592971/attachments/1925353/3187348/CERN_Site_Report_-_HEPiX_Autumn_2019.pdf

Заключение

- Расходы на информационную инфраструктуру растут и должны будут расти в мире, если исследовательская деятельность не прекратится.
- Техническая инфраструктура имеет тенденцию к концентрации в крупные вычислительные установки, которые обслуживают отдельные отрасли или географические регионы.
 - Внутри отдельных исследовательских групп рационально иметь микро-вычислительные установки.
- Растут затраты на специализированные устройства в составе детекторов/сенсоров в исследовательском инструментарии для отбора/фильтрации полезных данных. Уже сейчас число и производительность сенсоров/детекторов таковы, что порождаемый поток “сырой” информации переполнит любые каналы передачи и системы хранения данных.
- Граница между теоретической и экспериментальной наукой продолжает стираться, поскольку в процессе анализа нередко используются данные наблюдений/измерений и компьютерного моделирования совместно.

Литература 1

- <https://www.itelescope.net/>
- <https://rg.ru/2019/11/06/reg-cfo/sud-podtverdil-zakonnost-kamer-videoraspoznavaniia-lic.html>
- https://en.wikipedia.org/wiki/General_Data_Protection_Regulation
- https://www.imd.org/globalassets/wcc/docs/imd_world_digital_competitiveness_ranking_2018.pdf
- http://reports.weforum.org/delivering-digital-infrastructure/introduction-the-digital-infrastructure-imperative/?doing_wp_cron=1575992959.9462480545043945312500
- <https://archive.nytimes.com/www.nytimes.com/ref/membercenter/nytarchive.html>

Литература 2

- <https://www.vedomosti.ru/technology/articles/2019/11/26/817135-vlasti-sledit>
- <https://www.xsede.org/>
- <https://www.bnl.gov/compsci/SDCC/hpc/hpc1/index.php>
- <https://www.sdsc.edu/>
- <https://www.olcf.ornl.gov/>
- <https://computing.llnl.gov/>
- <https://www.bnl.gov/compsci/>
- Sample Size Planning for Classification Models
<https://arxiv.org/abs/1211.1323>

Литература - 3

- https://en.wikipedia.org/wiki/List_of_electronic_laboratory_notebook_software_packages
- <https://stats.stackexchange.com/questions/51490/how-large-a-training-set-is-needed>
- IEEE Transactions on Neural Networks and Learning Systems -
<https://ieeexplore.ieee.org/xpl/RecentIssue.jsp?punumber=5962385>