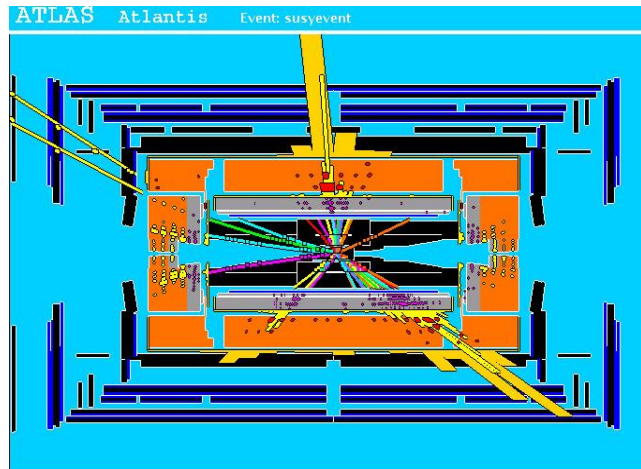


La Computación (¡ y el software !) en la era del LHC



Dr. José F. SALT CAIROLS
Profesor de Investigación CSIC
Instituto de Física Corpuscular
(Jose.Salt@ific.uv.es)



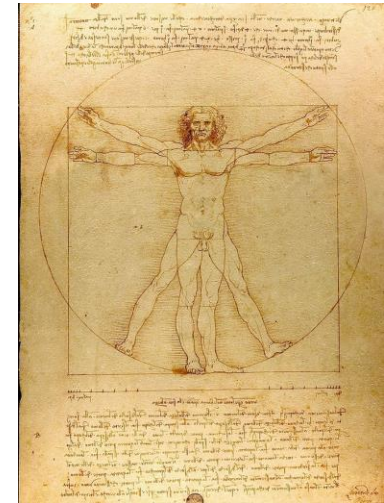
Contenido

- 1) **Introducción: Física y Computing**
=> Un poco de historia
- 2) **El Problema de la Computación y los Datos en el LHC**
- 3) **La Solución al problema: las Tecnologías GRID**
- 4) **El Modelo (Dinámico) de Computación de ATLAS**
- 5) **El Tier-2 Español de ATLAS**
- 6) **El GRID y el Higgs**
- 7) **Retos de Computación en el futuro.**
- 8) **Mensajes a retener/Conclusiones**

1- Introducción: Física y Computing

- **Físico Experimental (perfil amplio de un físico)**

- Hardware,
- constructor de detectores, R&D
- Calibración, Performance, Análisis de datos de los detectores
- Software
- Análisis de datos (Physics Analysis) y su interpretación
- Especialista en MC /Simulación
- Computing
-



- **La computación:** La tecnologías de la Información y las comunicaciones son pieza fundamental en muchas actividades. En la Ciencia su papel es esencial para la obtención de resultados



El origen (?)

- Física y Computación: confluyen en los años 30 y 40 del pasado siglo
- Un aspecto muy importante: la aparición de las Técnicas de Simulación (Monte Carlo) a partir de la concentración de científicos en Los Álamos cálculo de secciones eficaces, etc



Modelo Matemático de Implosión

John Von Neumann



Stanislaw Ulam

- Los físicos y matemáticos pioneros y visionarios:
 - Von Neuman, Ulam, Metropoli--- Fermi, Wigner, ...



Nicholas Metropolis



E. Fermi

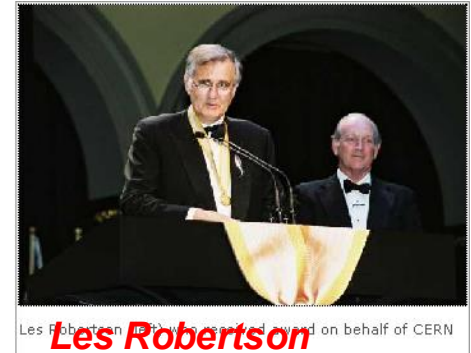


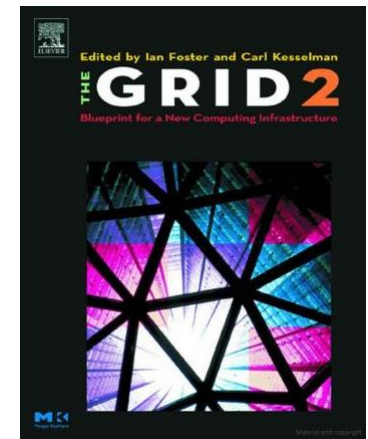
R.P. Feynmann

- Feynmann y la Informática Cuántica:
 - Charla en 1981, 'Simulating Physics with Computers' -> Propone el uso de fenómenos cuánticos para realizar cálculos computacionales (dada la naturaleza de complejidad de los cálculos)- Ordenador Cuántico

Evolución Computing Distribuido

- **Agregación de recursos:**
 - Sistema SHIFT (“granjas” RISC Unix, L.Robertson, CERN)
 - Clusters de PCs (Sistemas Beowulf, “fábricas” de PC...)
- Recursos **distribuidos compartidos:**
 - Sistema Condor (M.Livny)
 - gestión de tiempo inactivo en sistemas Linux de la red local
 - Red Entropía, Programa SETI@home
 - Sistemas “Peer to Peer”
- Aparición de **la Web:**
 - 1989: primera propuesta, CERN, Tim Berners-Lee y Robert Cailliau
 - Primeros servidores web en laboratorios de Física europeos
 - 1991: un prototipo del sistema WWW suministrado para la comunidad de HEP





La 'Biblia' del GRID



'Granjas' de ordenadores



- **1995: Supercomputing '95**
 - Experiencia I-WAY: 17 centros USA conectados a 155Mbps
 - Primeras iniciativas Grid:
 - NASA Information Power Grid
 - Iniciativa de la NSF con los centros NCSA y SDC
 - Advanced Strategic Computing Initiative (DoE)
- **La era de las los clusters/granjas/fabricas de PCs**
 - Nodos con CPUs (duales)
 - Disco local + acceso a servidores (NFS-NAS, AFS)
 - Limitación:
 - Interconexión de red: Gigabit casi popular, pero Latencia baja requiere Myrinet o similar, solución mas costosa
 - Perfectos para HTC (High Throughput Computing)
 - Las aplicaciones HPC (High Performance Computing) necesitan adaptarse:
 - La memoria no está compartida
 - Las herramientas: PVM, MPI
 - Instalación y control: funciona bien para cientos de ordenadores “homogeneos”

- ... Y además de forma distribuida mediante la conexión de clusters
 - Gracias a la interconexión de la Red
 - Llegar a obtener la sensación de que trabajamos con un Superordenador formado por varios clusters de ordenadores (Metacomputación)



2- El Problema de la Computación y los Datos en el LHC

Almacenamiento-

Ratio de registro de datos 0.1 - 1 GBytes/sec

Acumulando a 10-15 PetaBytes/año

15 PetaBytes de disco

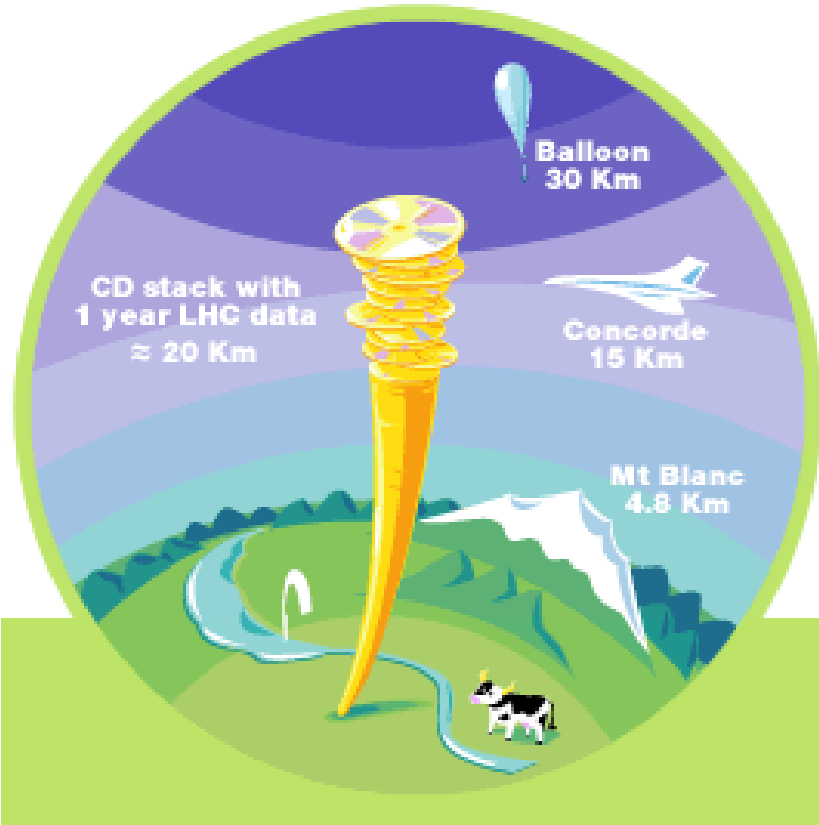
Procesamiento

100,000 de los PC's más rápidos actuales



- La infraestructura computacional del CERN **NO** era suficiente para procesar todos los datos

- Se necesitaba un sistema 'robusto': que haya redundancia, sin 'únicos puntos de fallo'



Los 4 experimentos del LHC:
Europa: 267 institutos, 4603 usuarios
Resto mundo: 208 institutos, 1632 usuarios

Tecnologías Computación GRID



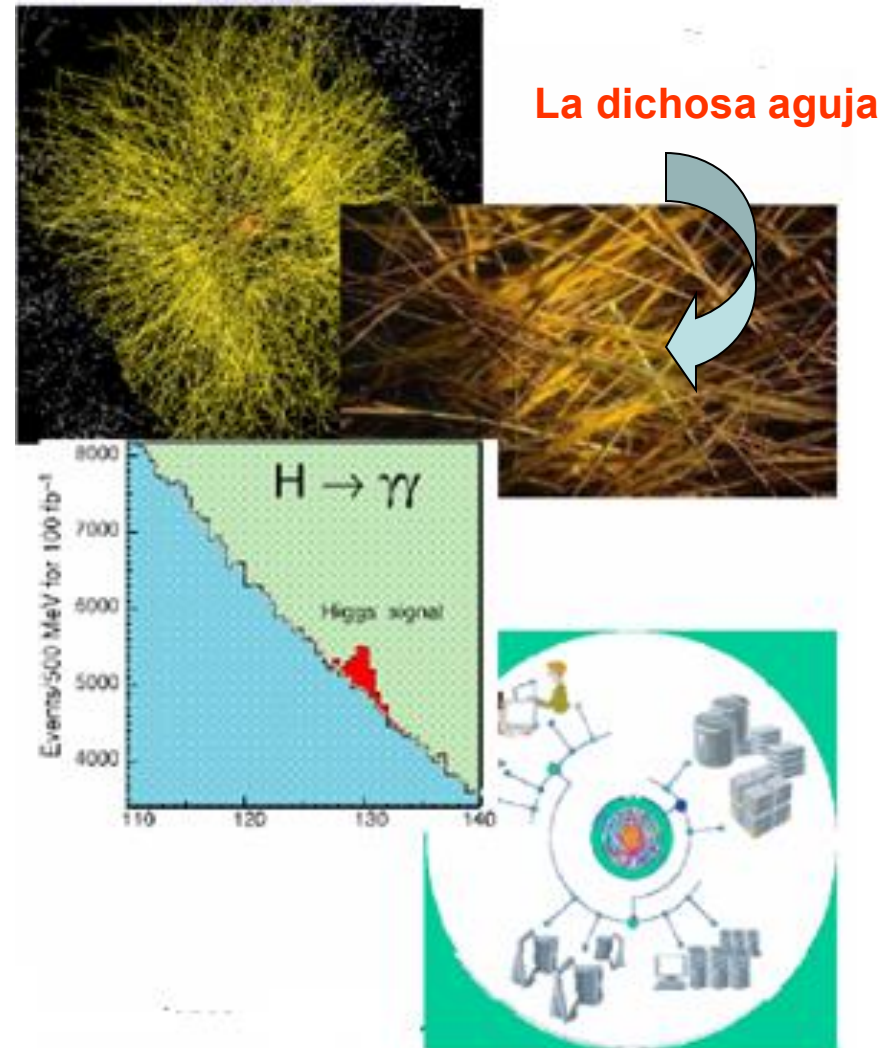
3.- La Solución al Problema: las Tecnologías GRID

¿Qué es la Computación GRID?

- **Definición:** la computación GRID es una evolución de la computación distribuida: su base está en tecnologías que permiten a las organizaciones compartir recursos informáticos para hacer frente a diferentes tipos de necesidades . En general, dichos recursos están dispersos geográficamente y conectados por Internet pero dicha red proporciona la impresión de estar trabajando con un único sistema informático virtual
- **intenta resolver problemas actuales de la Sociedad de la Información:**
 - Acceso rápido a bases de datos/almacenamiento...
 - Proporcionar su procesamiento y análisis utilizando potencia de cálculo distribuida y potentes facilidades de visualización...
 - Mediante la utilización de la red (Internet)

Las Tecnologías de la Información nos ayudan a 'buscar una aguja en un pajar'

- La Señal es extremadamente baja (**1 sobre 10 billones de colisiones**)
- Volumen de datos:
 - (Frecuencia alta) *(gran número de canales) * (4 experimentos) =
15 PB de nuevos datos/año
- Potencia de cálculo:
 - (Complejidad de sucesos) * (número de sucesos) *(miles de usuarios)=
100 K de los CPUs más rápidos actuales
45 PB de almacenamiento de disco (Simulación +Análisis)
- Análisis y financiación mundial:
 - Financiación local de recursos de computación en los principales países y regiones



Simil del GRID Computacional con la Red Eléctrica

Red Eléctrica

- Plantas de producción de electricidad



- Distribución jerárquica del flujo eléctrico



Subestación Eléctrica

- Tendidos eléctricos



- Usuario final: acceso a las prestaciones de la electricidad



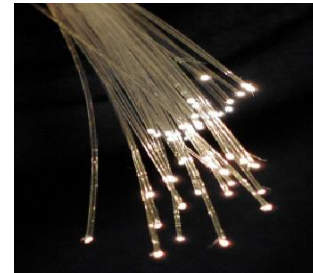
• GRID Computacional

- Grandes centros de almacenamiento de datos/potencia de cálculo
(Supercomputador/GranCentro de Cálculo)

- Distribución jerárquica
(Centros de cálculo medianos)

- Red Internet

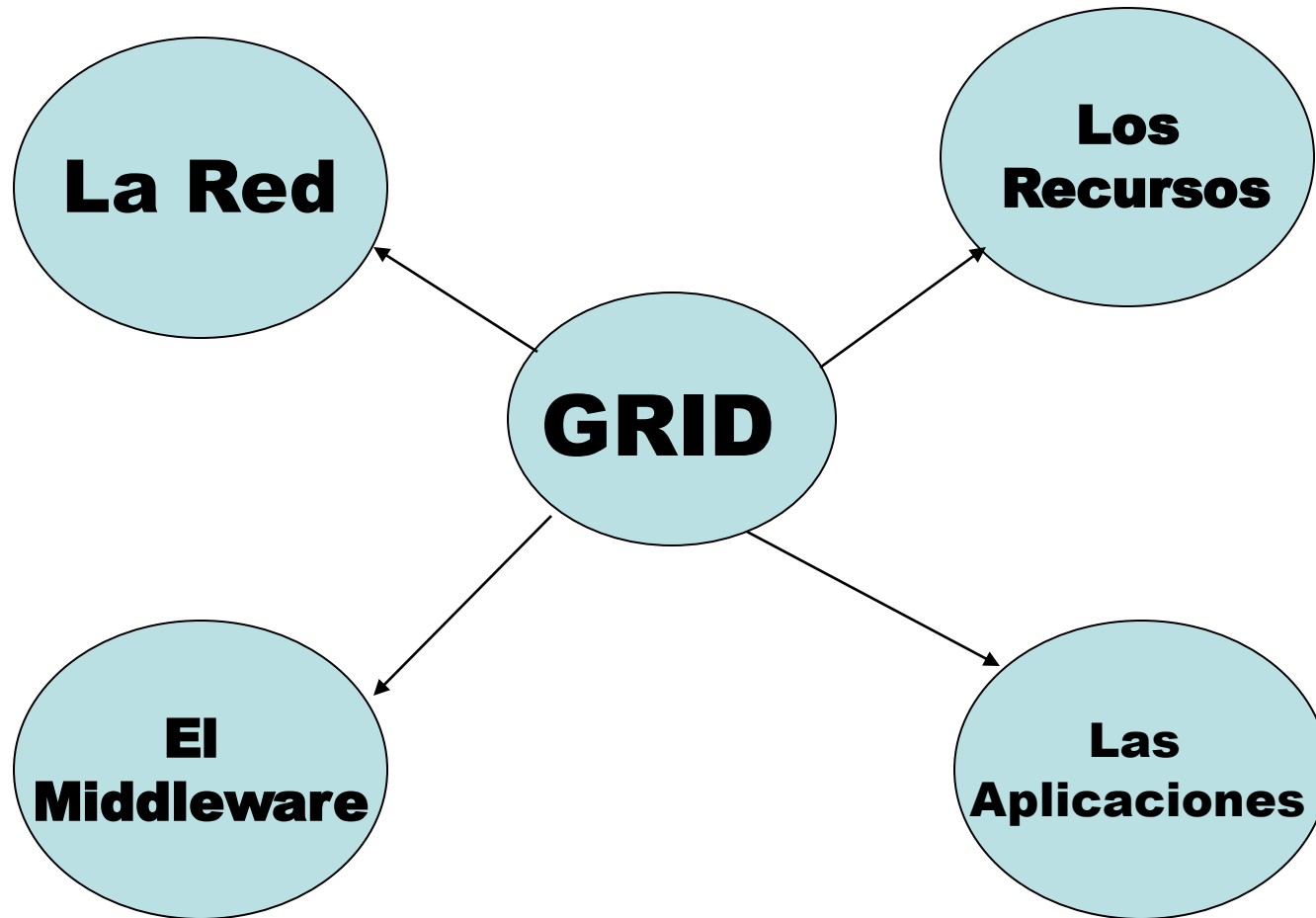
← Ramo de fibras ópticas



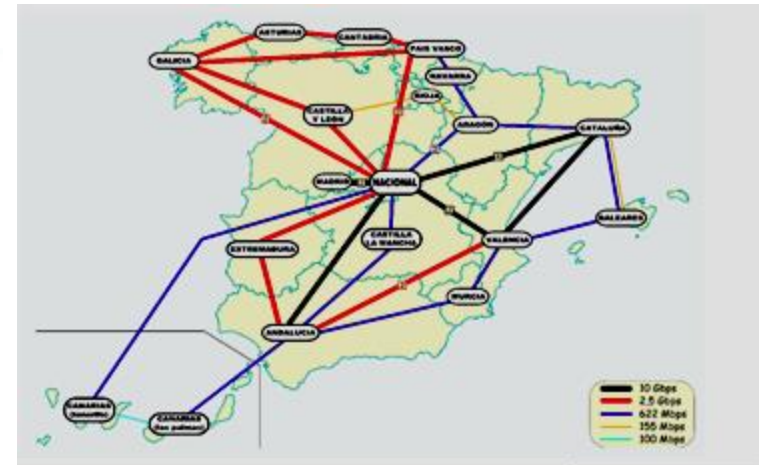
- Usuario Final: acceso a las prestaciones informáticas



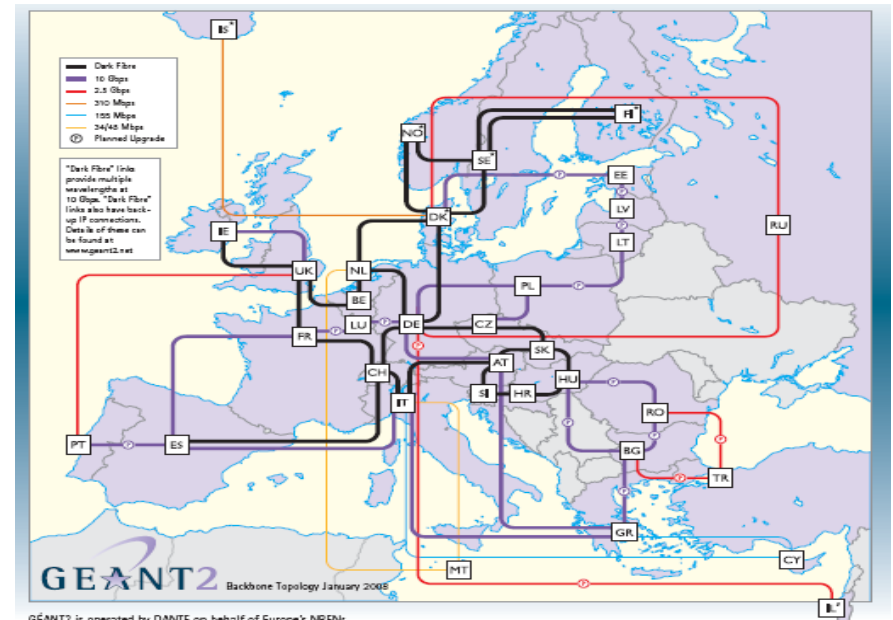
Ingredientes fundamentales del GRID



La Red: las autopistas de la Información



Mapa de la Infraestructura de la Red en España (RedIris)



La RED Académica y de Investigación Pan-europea

- El GRID no podría desarrollarse sin una Red apropiada.
- Responsable de asegurar los recursos que forman el GRID
- Comunicación
 - Protocolos de Internet: IP, DNS, routing, etc.
- Gran esfuerzo a nivel Europeo (DANTE y NRENs) y Español (RedIris)
- Constituyen **las Autopistas de la Información**

Los recursos: Las Infraestructuras

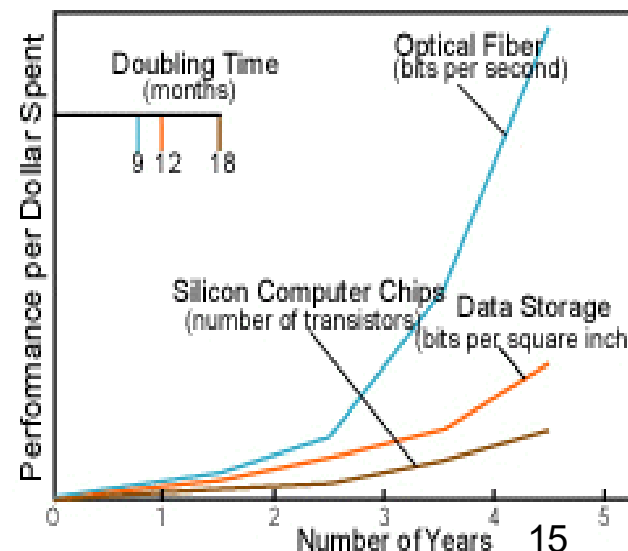
Evolución tecnológica

- **Hardware**
 - Procesadores...cumplen la ley de Moore: se **duplica cada 16-18 meses !**
 - Nodos de computación:
 - PCs y portátiles, Estaciones, Servidores, Clusters, Blades
 - Almacenamiento: **se duplica cada 12 meses !!**
- **Mejora de la Red:**
 - la capacidad de la red **se duplica cada 9 meses !!**
 - 10 GB ethernet
- **Coste tiene en cuenta muchos factores:**
 - Espacio, alimentación eléctrica, mantenimiento

- **Tendencia al 'Commodity Computing' :** adquisición de sistemas de diferentes compañías (vendedores), incorporan componentes basados en standards 'abiertos'



Gordon Moore

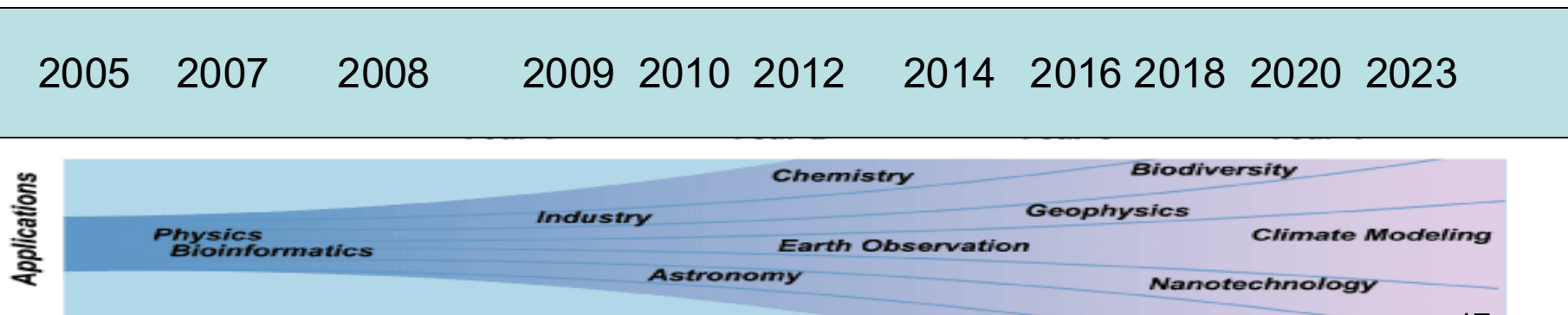


¿Qué es el Middleware?

- **es un software de conectividad que ofrece un conjunto de servicios que hacen posible el funcionamiento de aplicaciones distribuidas sobre plataformas heterogeneas.**
- **nos abstrae de la complejidad y heterogeneidad de las redes de comunicaciones subyacentes, de los sistemas operativos y lenguajes de programación,**
- **Finalidad: ‘virtualizar’ los recursos de computación’**

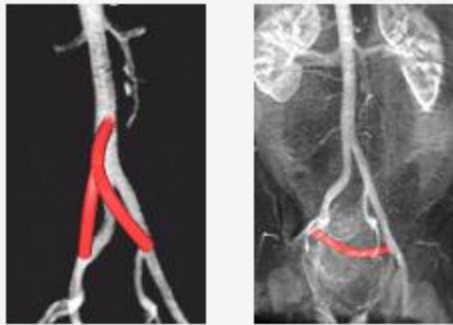
Las Aplicaciones

- Algunas disciplinas científicas se han organizado durante las décadas pasadas en grandes colaboraciones científicas que agrupan una gran cantidad de científicos al ser la única manera de alcanzar logros de alto nivel que no sería viable con grupos de trabajo reducidos
- resolución de retos tecnológicos importantes (ejemplo: experimentos del LEP, Proyecto Genoma, etc)
- En aspectos computacionales y tecnologías de la información, se ha evolucionado y ha desembocado en el paradigma GRID.

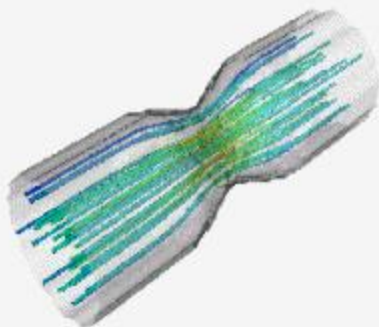


El GRID aplicado a enfermedades vasculares

- Dos procedimientos: stent y bypass
- El cirujano se ayuda con imágenes 3D obtenidas con scanners CT y MRI
- Las simulaciones de una posible intervención ayuda a la toma de decisión por parte del médico

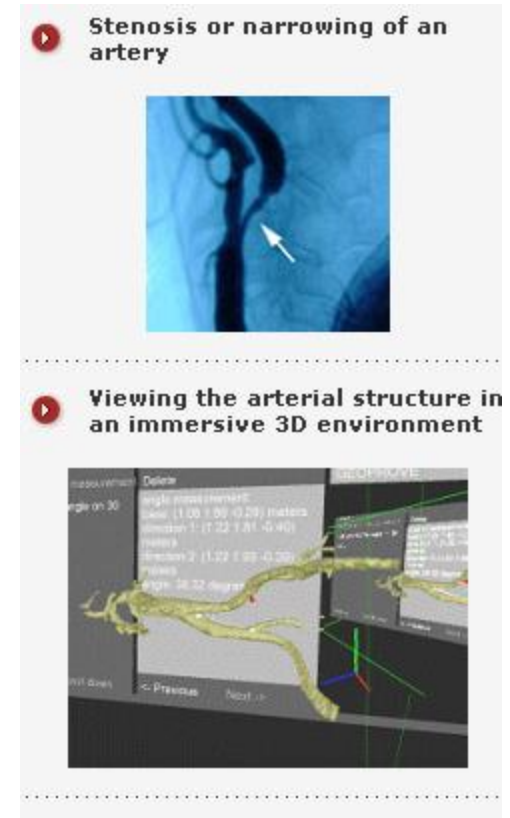


▶ Simulated flows



El GRID permite lanzar los cálculos pertinentes. Se utilizan los recursos de los centros que participan en el proyecto (CROSSGRID)

Investigación coordinada por la Universidad de Amsterdam y con la colaboración del CSIC y de todos los centros de CROSSGRID



4.- El Modelo Dinámico de Computación de ATLAS



- El proyecto de Computación GRID para el LHC (WLCG) comenzó en 2002:
 - Fase I (2002-2005): test y desarrollo, construcción y prototipo de servicios
 - Fase II (2006-2010): despliegue inicial de los servicios GRID
 - Fase III (2011-2022): consolidación de servicios, ampliación a otros paradigmas de computación
- Propósito: Suministrar los recursos de computación necesarios para procesar y analizar los datos recolectados por los experimentos del LHC

Enormes volúmenes de datos y capacidad de computación

- 15 PB/año (Raw) -> 50 PB/año (overall)
- Vida Util: 10-15 años -> escala del Exabyte

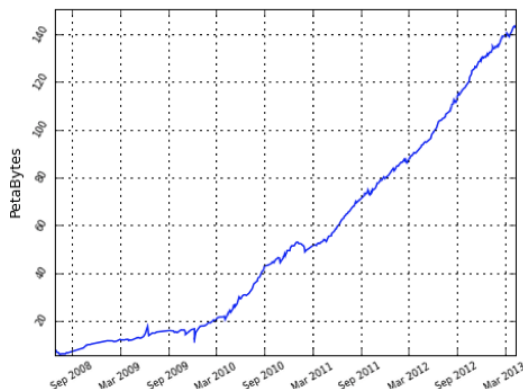
Reto de la Escalabilidad

- Resulta evidente que se necesita una infraestructura distribuida

Logros:

- CMS: ha transferido 100-200 TB/día por el GRID en los últimos dos años
- ATLAS: se han acumulado 140 PB durante el procesamiento de datos del Run I

Total GRID space usage according to DQ2



Organización de centros de recursos en ATLAS: Tiers

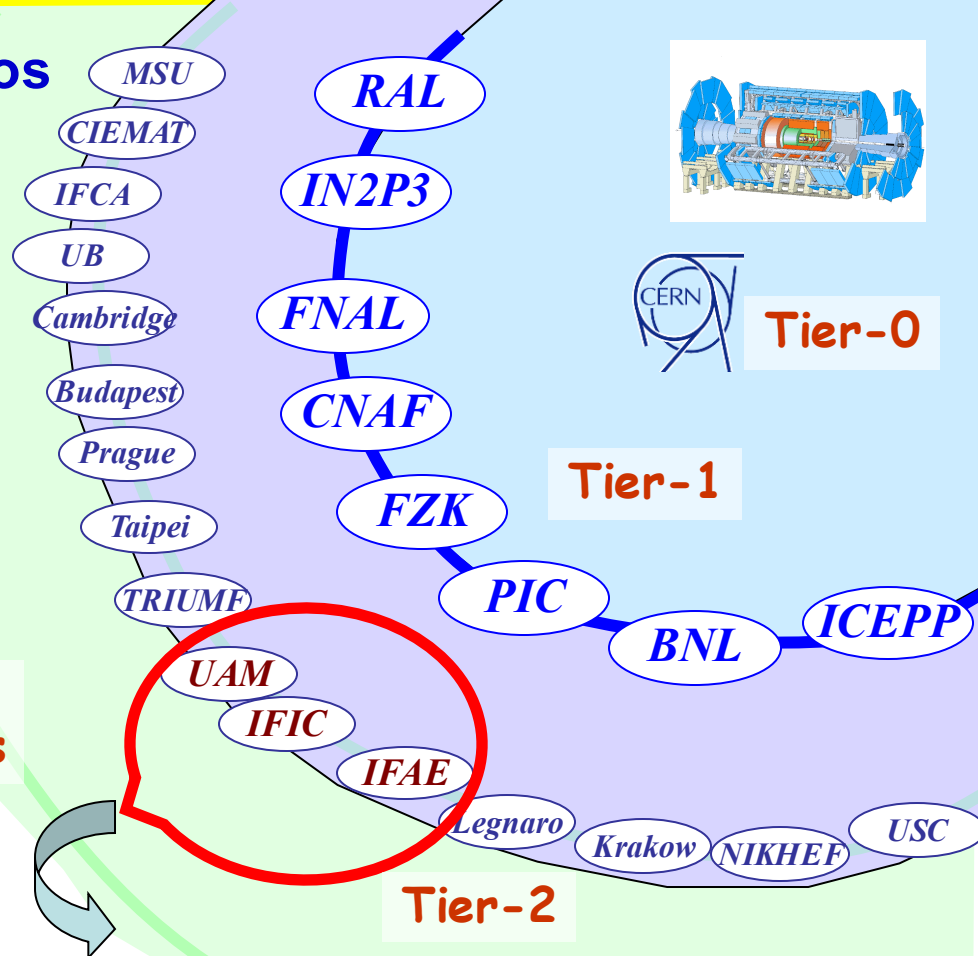
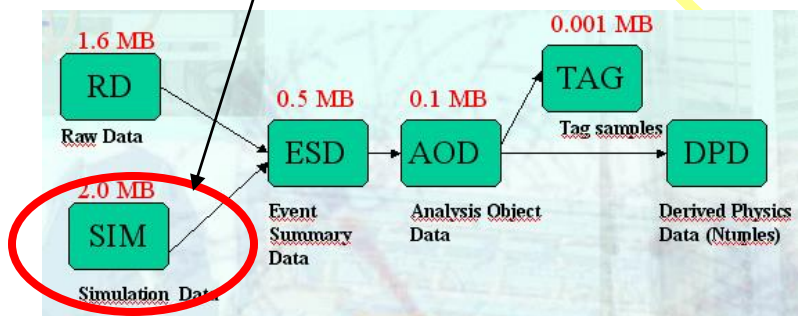


desktops
portables

Tier-3

small
centres

- Los datos sufren transformaciones encaminadas a la reducción en tamaño y la extracción de la información relevante
- Con el fin de prepararnos para los análisis con Datos Reales y testear el sistema, trabajamos con Datos Simulados (Monte Carlo)



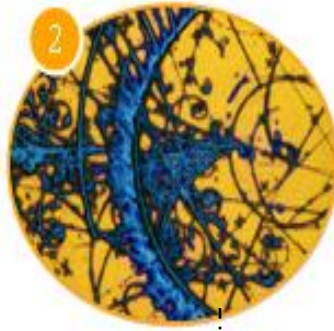
asegurar un crecimiento sostenible de la infraestructura Tier-2 entre estos centros y su operación de forma estable

Proceso de transformación de los datos



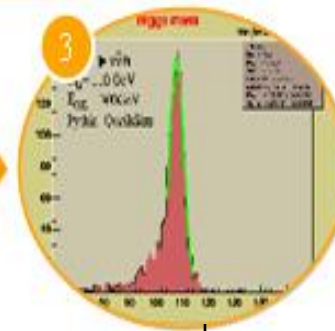
Datos de ATLAS 'brutos'

Proceden directamente del experimento (después del "trigger") ubicados en el CERN, serán procesados y pasados a los Tier-1 (1.6MB/colisión)



Datos procesados

Proporcionan información clara de trazas y valores de energía registrada por los detectores del experimento ATLAS. (500KB/colisión)

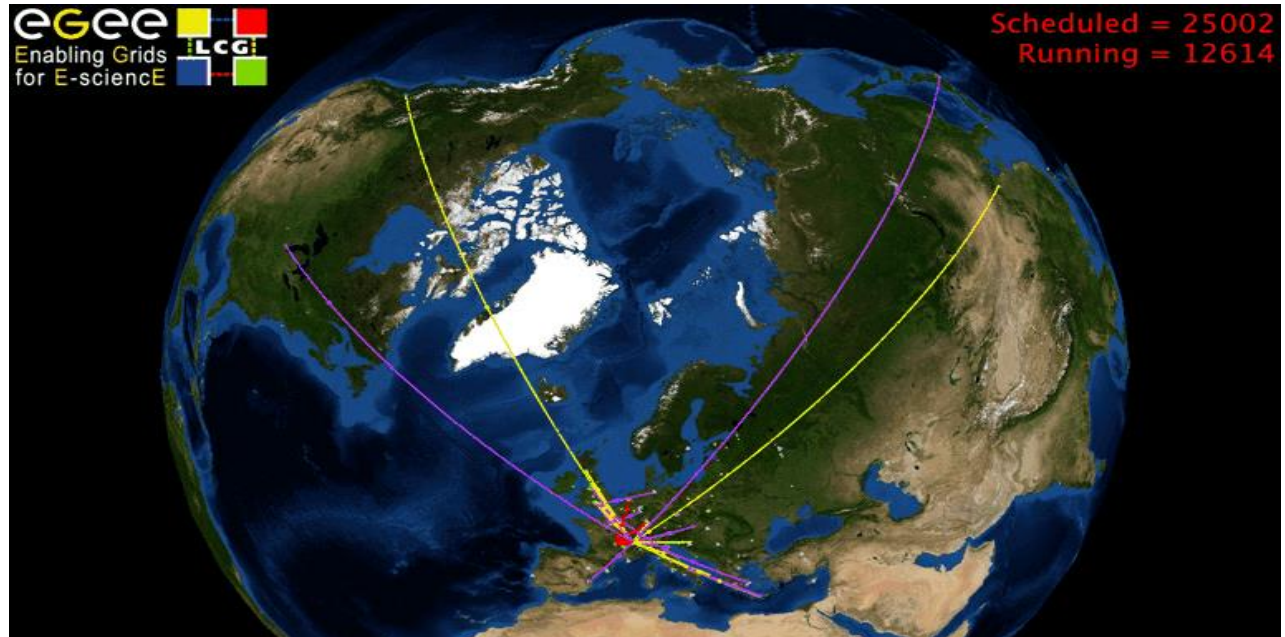


Datos listos para analizar

Provenientes de los datos procesados, servirán para descubrir nuevas partículas por los científicos de ATLAS. (100KB/colisión).

- Los datos sufren transformaciones encaminadas a la reducción en tamaño y la extracción de la información relevante



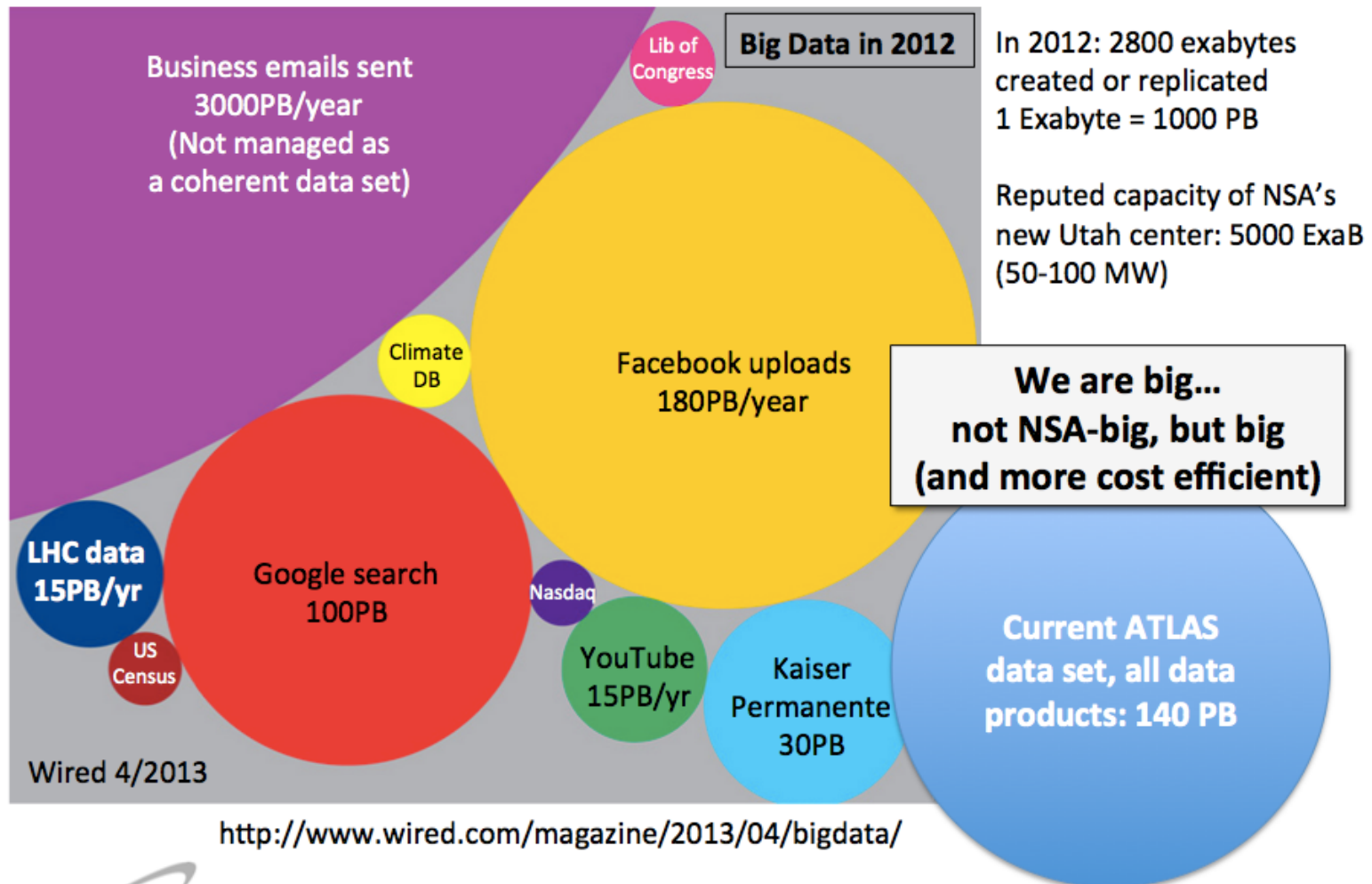


15:

**El GRID: “Un despliegue colaborativo donde no se pone nunca el Sol”
(‘round the clock’ : 24 h / 24 h)**

Data Management

Where is HEP in Big Data Terms?



5.- El Tier-2 Español para ATLAS

Proyecto Coordinado del Plan Nacional de FAE de 3 centros:



IFAE (25%)

UAM (15%)

IFIC (60%)

4% del total

Tier-2

Equipamiento:

IFIC

CPU = 41.298. HEPSpec06

Disk = 4.032 TB

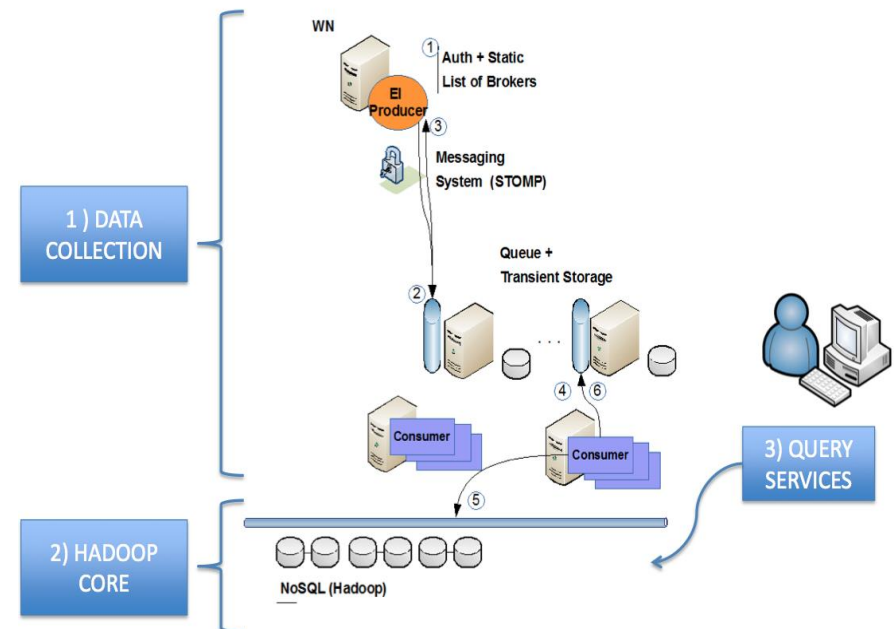
Recursos Humanos:

10 FTE

Julio 2023

Desarrollo del Event Index de ATLAS

- Desarrollo y despliegue de un catálogo completo de todos los sucesos de física (Real Data y MC) para todas las etapas de procesamiento
- Uso de nuevas tecnologías de Big Data y NoSQL (Hadoop)
- Alta participación del grupo IFIC:
 - Responsable of Data Collection Task.
 - Responsabilidad en la tarea de Data Collection
 - Ayuda en las verificaciones de consistencia en la producción global de datos



Evolución del Crecimiento de los recursos de TIER-2 para el experimento ATLAS



Primera Granja de PC's



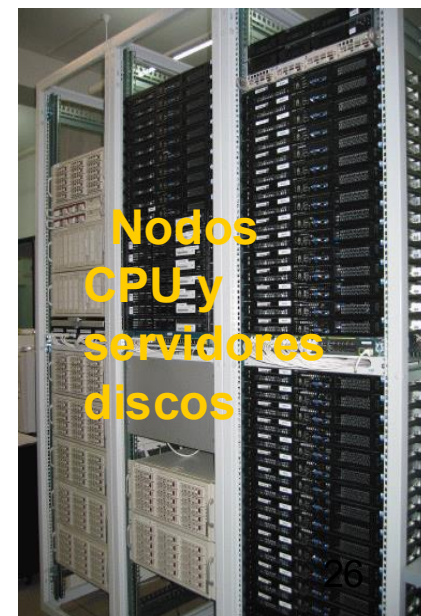
Actualización
Infraestructura
Computación del IFIC:
GLUON



Fuerte Crecimiento en recursos



Robot de
Almacenamiento
en Cintas

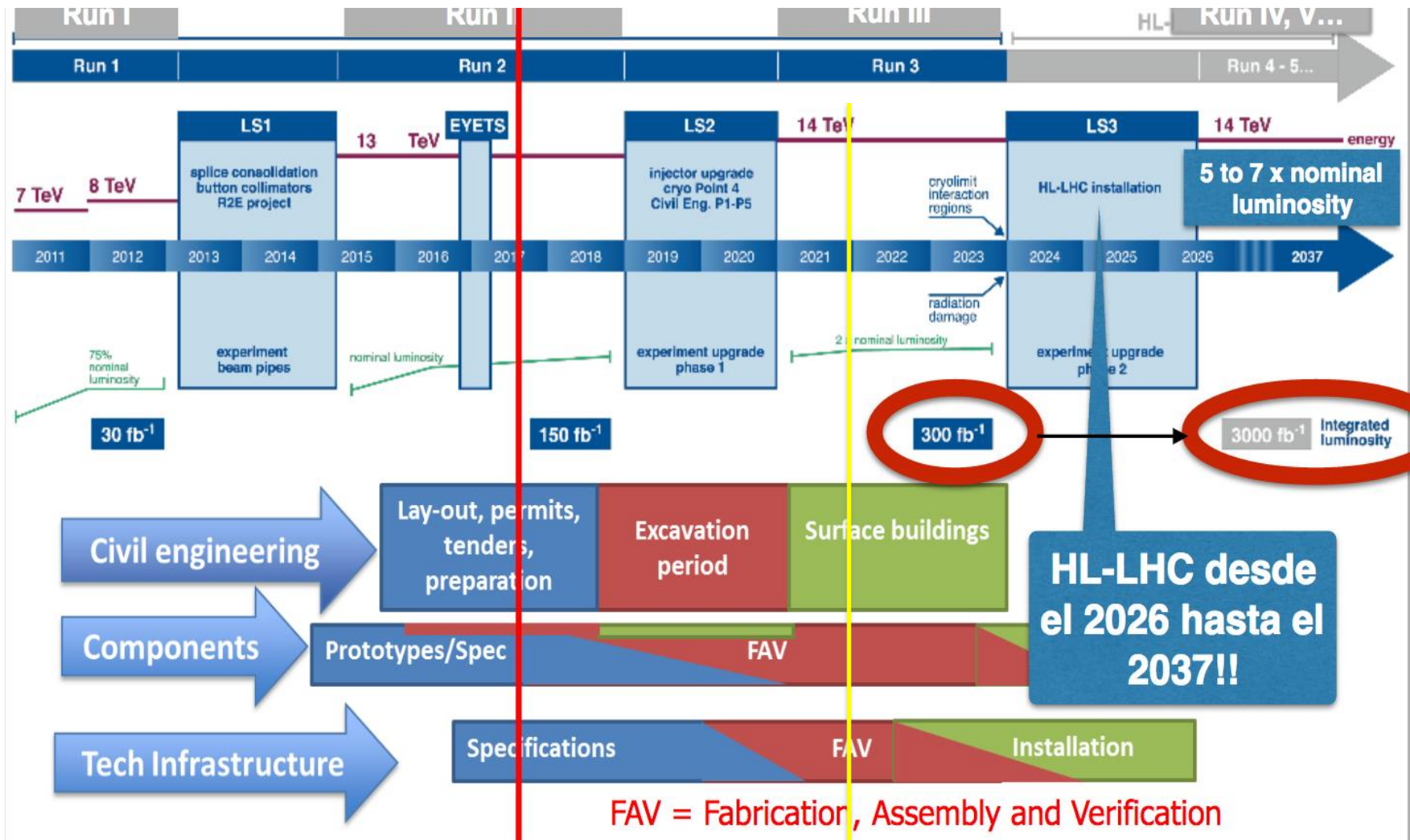


Nodos
CPU y
servidores
discos

6- El GRID y el Higgs

- Los resultados de Física son comparaciones de observables con sus predicciones teóricas/modelo
- Hay un proceso complejo entre la toma de datos y los resultados de física: Trigger/adquisición de datos, reconstrucción , análisis y simulación, grupos de análisis, publicación de resultados.
- Se necesitan algoritmos sofisticados para reconstrucción y análisis
- Uno de los factores más esenciales es disponer de una buena simulación detallada
- Un solo suceso tiene un determinado tamaño y 'simple', pero la escala de computación es enorme
- El software y el computing deben funcionar muy bien si se quiere lograr los resultados de primer nivel en el LHC

7- Retos de Computación en el futuro



Estamos aquí:
Julio 2023

El escenario actual: un 'ecosistema' de Computación

Estamos en un nuevo escenario debido a:

- las muy buenas prestaciones de **la red** en términos de latencia y ancho de banda
- la **mejora de la ratio calidad /coste** en los equipos de ordenadores
- Los **supercomputers (HPC)** han seguido una tendencia de crecimiento continuo durante las dos pasadas décadas
- **Cloud Computing** ha parecido hace unos 8-9 años como otra posible solución para resolver problemas actuales de computación
- **Sostenibilidad de la Infraestructura:** mayor rendimiento con menos energía, huella de Carbono, green computing

Computer cluster



Cloud Computing

GRID Computing



Supercomputers

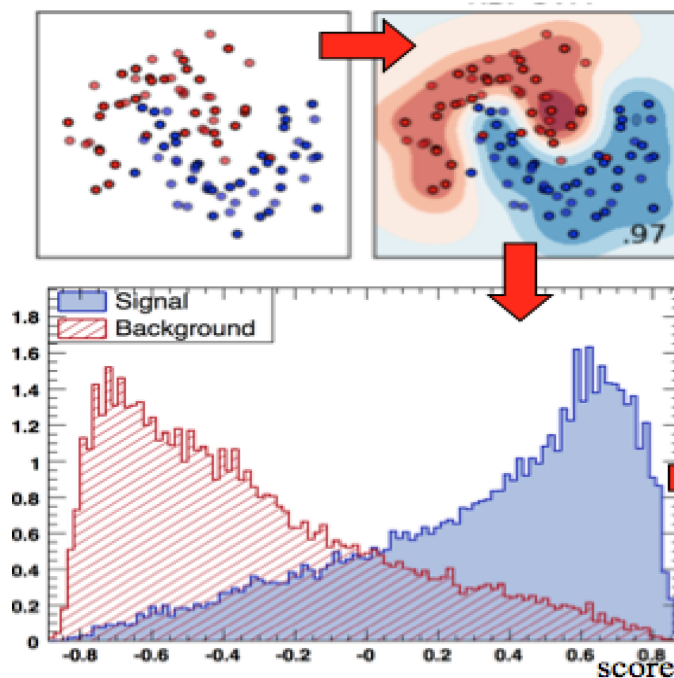
Machine Learning

- Utilización del Machine Learning en diferentes aspectos:

- ML en simulación
- ML en reconstrucción
- ML en análisis



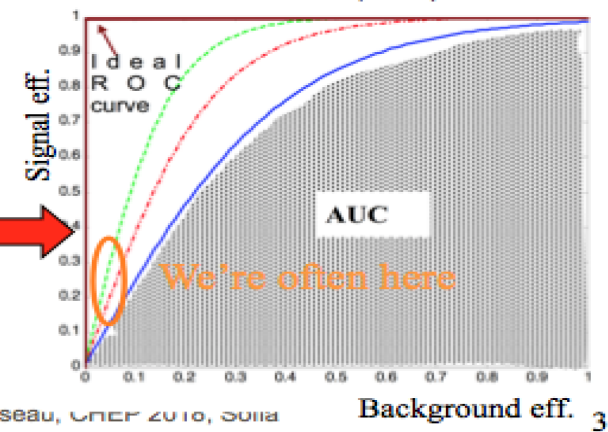
Clasificación Señal- Background



Train on Signal and Background Monte-Carlo
→ learn the separation between S and B distribution
Apply on test sample
Apply on data

Note: instead of classifying 0 or 1, can regress !

AUC : Area Under the (ROC) Curve



Salidas Profesionales:

A destacar dos (pero hay más!):

•En Investigación en Física Experimental de Partículas, Nuclear y Astropartículas (estáis viendo qué hacemos):

- trabajando en líneas de investigación que ya estáis conociendo en el IFIC Summer 2023

•Data Scientist:

- Profesión con muy buenas perspectivas y con una demanda cada vez mayor por parte de empresas y proyectos**
- Big Data: Obtención de conocimiento a partir de los datos**
 - El experimento ATLAS y otros experimentos (AGATA, DUNE, por ejemplo) son claros ejemplos de Big Data**
- Data Mining, Machine learning, Inteligencia Artificial, etc**

8.- Mensajes a retener/Conclusiones

- **La Física subatómica (HEP & Nuclear & Astroparticle) y Computación están estrechamente relacionadas**
- **La Computación en la era del LHC se ha desarrollado para resolver problemas originados por las características de los experimentos. La solución inicial ha sido el GRID**
- **El Computing es como un detector más : colaboración de sites, upgrades, organización de servicios y de 'shifts', R&D , etc**
- **El I + D + i asociado a estas tecnologías GRID está dando lugar a una 'segunda revolución' del mismo nivel que el que supuso la web hace 20 años**
- **Esto está generando un Modelo Dinámico de Computación**
- **Aplicaciones en otras disciplinas científicas. Ya se están obteniendo retornos en Biomedicina, Sector Energético, Química, Biotecnología, etc**
- **El GRID y otras contribuciones de las Tecnologías de la Información y de la comunicación han dado lugar a la e-Ciencia**
- **En Investigación es necesario ' Data Scientists' y son también muy apreciados en las empresas en el sector privado y público.**

- ChatGPT
- Operación de las infraestructuras con IA

- GRACIAS por vuestra atención
- PREGUNTAS / COMENTARIOS

BACKUP SLIDES

Grupo de Computación GRID y e-Ciencia (G2C2E)

Group of GRID & e-Science

Personal Permanente:

Álvaro Fernández Casani	Titulado Superior Informática	CSIC
Santiago González de la Hoz	Catedrático Universidad	UV
José Francisco Salt Cairols	Profesor de Investigación	CSIC
Francisco Javier Sánchez Martínez	Titulado Superior Informática	CSIC

Personal Contratado o Vinculado:

Carlos García Montoro	Contratado por Proyecto	CSIC
Javier Aparisi Pozo	Contratado Pre-Doc S8A	CSIC
Esteban Fullana Torregrosa	Contratado Prometeo (postdoc)	CSIC
Miguel Villaplana Torregrosa	Investigador CIDEAGENT	UV
Victoria Sánchez Sebastián	Contrato Pre-Doc FPU	CSIC

Farida Fassi	Profesora Titular Universidad
	De Rabat

Colaboradores Externos/Colaboradores Visitantes:

Julio Lozano Bahilo
Gabriel Amorós Vicente
Mohammed Kaci