

Sistema de análisis distribuido de datos procedentes de ATLAS utilizando el GRID

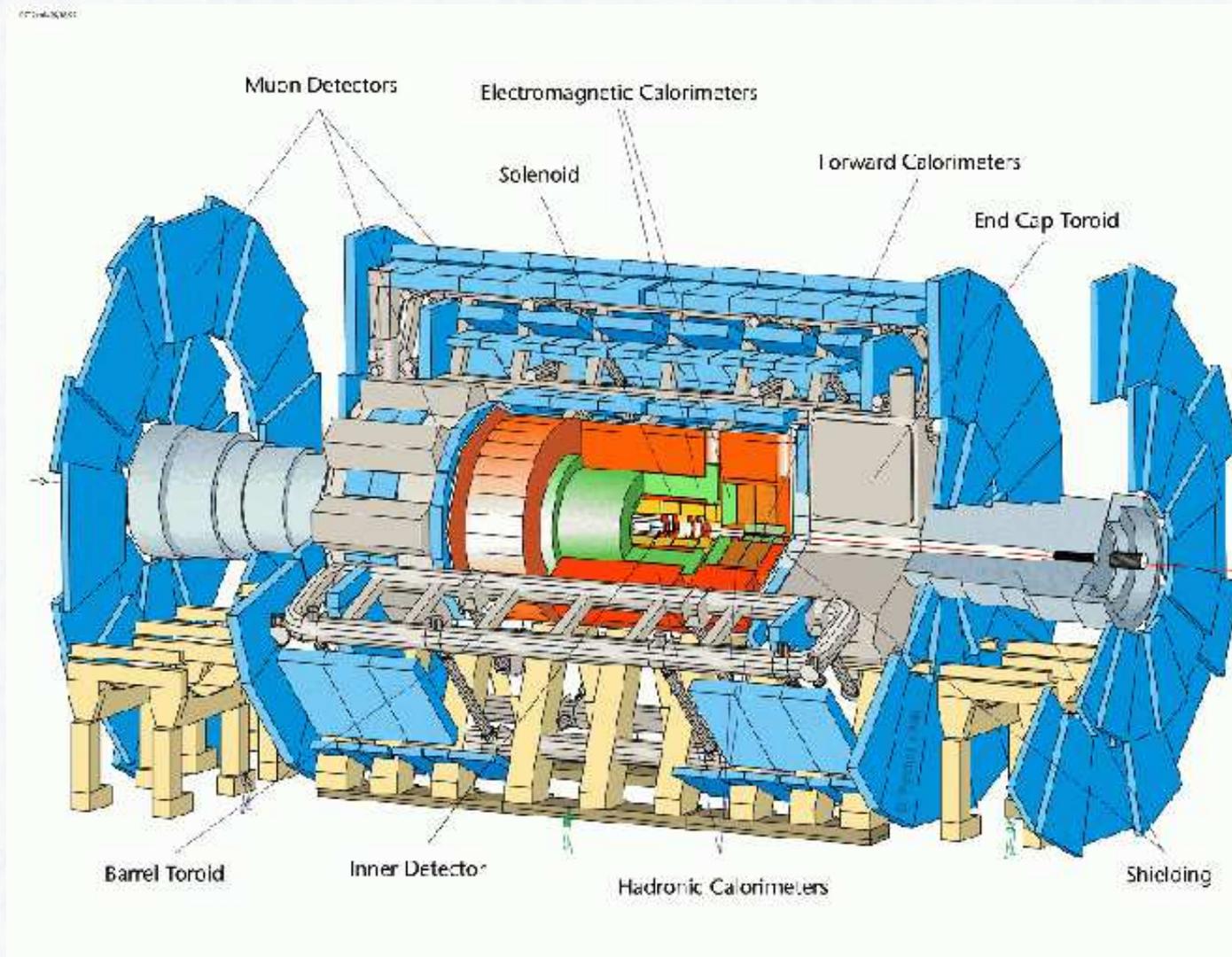
Ourense, 12-16 de Septiembre de 2005

F. Fassi, J. Lozano, L. March, J. Salt, S. González,
J. Sánchez, A. Fernández, M. Kaci y M.D. Jordán
Instituto de Física Corpuscular
(CSIC/UV)

Contenido

- ✓ Introducción
- ✓ El Sistema de Análisis Distribuido de ATLAS (ADA)
 - AJDL (Abstract Job Description Language)
 - Arquitectura del sistema
 - ❖ El entorno de usuario
 - ❖ Servicio de catálogos
 - ❖ Servicios de análisis
- ✓ Conclusiones

Introducción



El Sistema de Análisis Distribuido de ATLAS

- ✓ Proporciona un sistema de fácil acceso y con una funcionalidad que permite realizar el análisis en un entorno donde los usuarios, los datos y el procesamiento están distribuidos geográficamente.
- ✓ Toma y combina las herramientas de otros proyectos para completar su arquitectura
- ✓ Se basa en una colección de Servicios WEB donde el usuario interacciona con el sistema mediante una interfaz que implementa los componentes y que permite:
 - ocultar los detalles y la complejidad del middleware que utiliza ADA
 - facilitar el empleo del sistema para:
 - el envío de los trabajos y su monitorización
 - el acceso a los catálogos, a los ficheros y al software necesario para las transformaciones.

AJDL

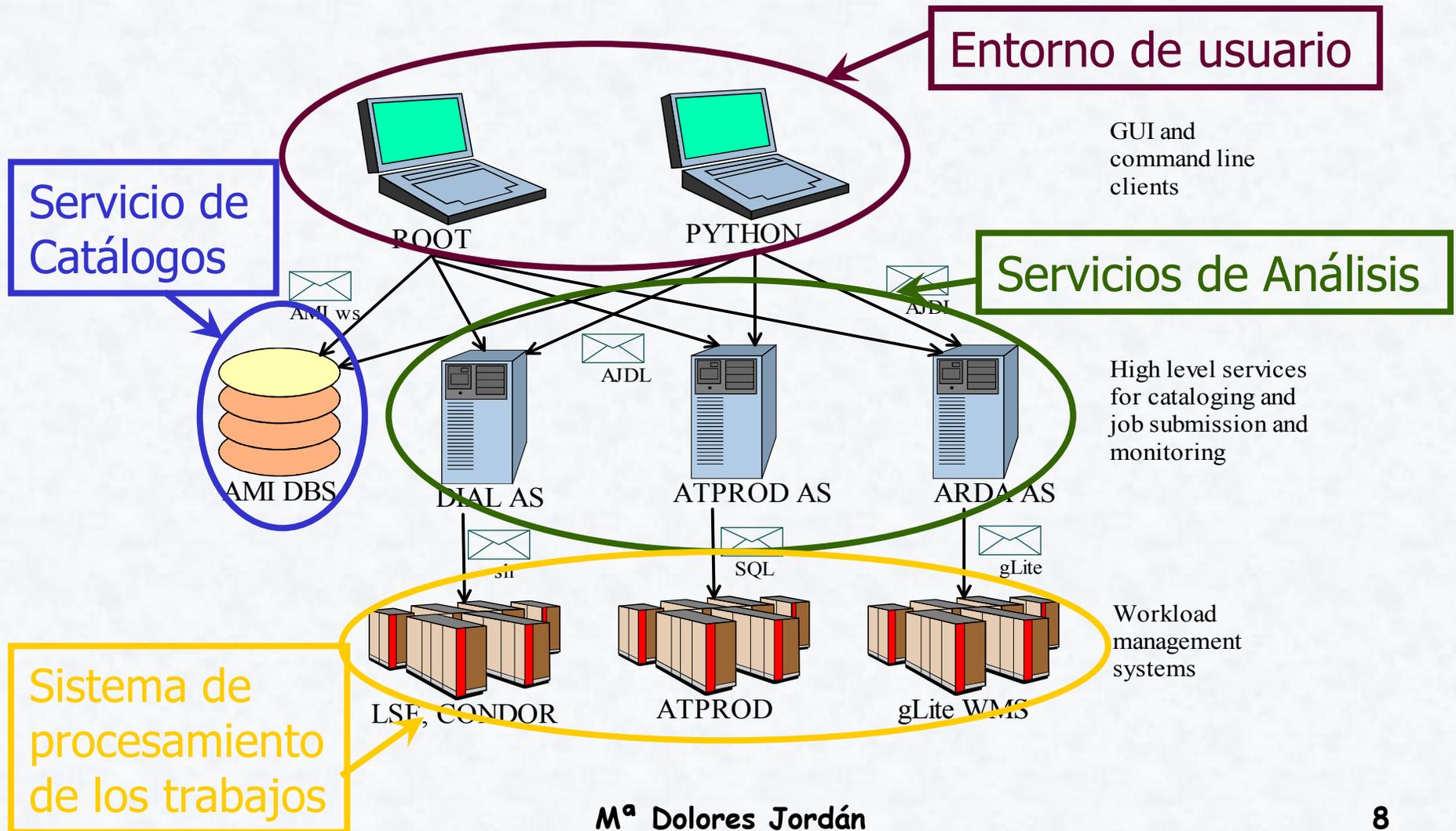
Abstract Job Description Language

- ✓ El usuario, desde su entorno, puede comunicarse con los distintos servicios que proporciona el sistema a través del uso de un lenguaje abstracto de descripción → AJDL
 - ✓ El AJDL está implementado en clases de C++
 - o Dataset
 - o Transformación
 - o Trabajo
- Definen los componentes del sistema

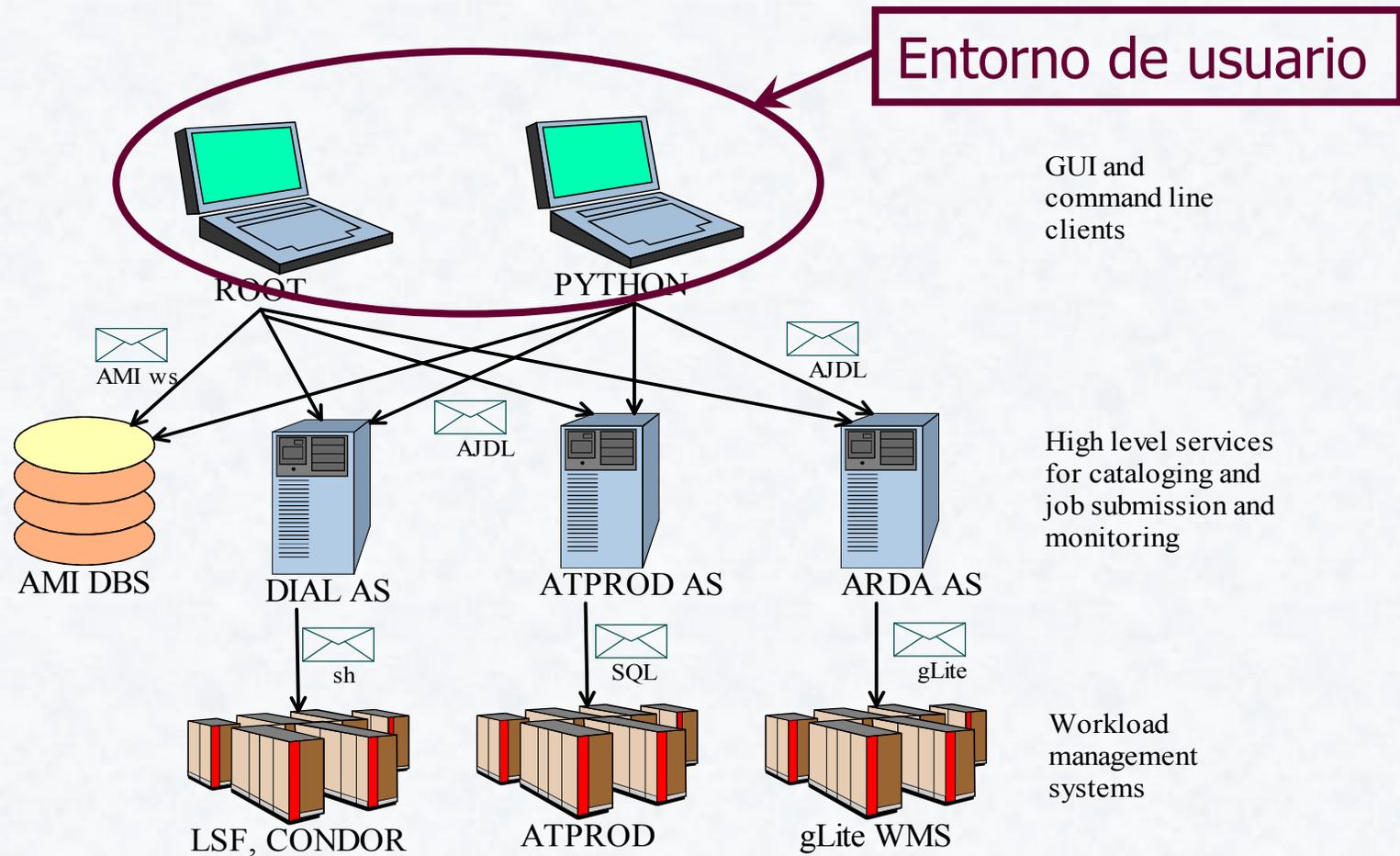
Componentes de AJDL

- Dataset
 - ❖ describe los datos
 - ❖ proporciona al usuario información para:
 - ❑ localizar datos
 - ❑ conocer su tipo, y el número de sucesos que contiene
- Transformación
 - ❖ permite la manipulación y la explotación de los datasets
 - ❖ describe una operación que actúa sobre un dataset produciendo un nuevo dataset
 - ❖ Está dividida en:
 - ❑ **Aplicación** → especifica los ejecutables y las librerías compartidas
 - ❑ **Tarea** → especifica la configuración del usuario para una aplicación
- Trabajo:
 - ❖ Ejemplo particular de transformación que actúa sobre un dataset

Arquitectura de ADA



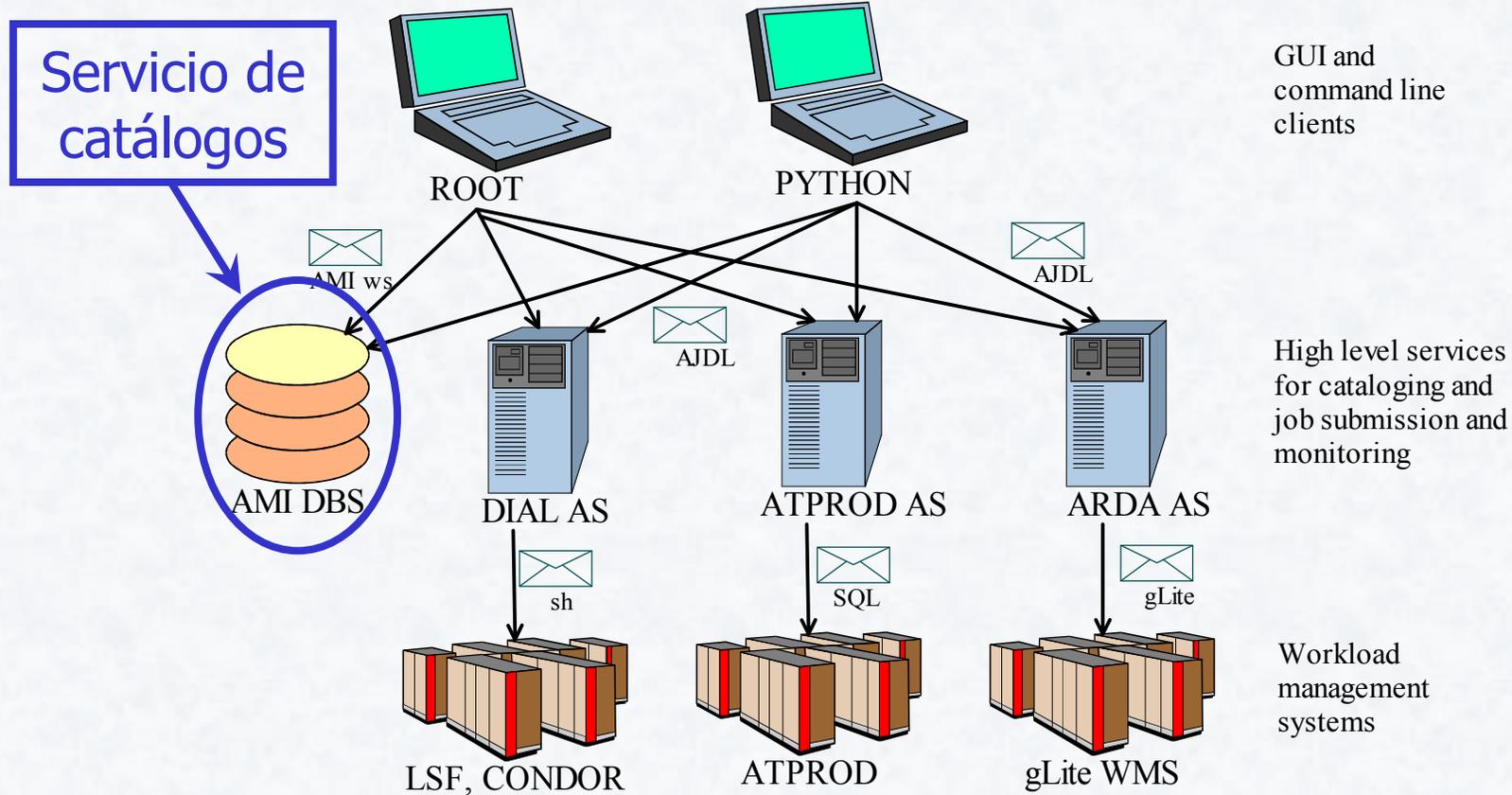
Entorno de usuario



Entorno de usuario

- ✓ Proporciona al usuario los medios para comunicarse con los servicios de ADA
- ✓ Permite:
 - acceder fácilmente a los datos
 - monitorizar el progreso de los trabajos
 - obtener un resultado parcial o completo de los trabajos
- ✓ Facilita al usuario los recursos para examinar y crear los datasets, transformaciones y trabajos.
- ✓ Se accede a la implementación de los componentes de AJDL bajo:
 - ROOT implementado en C++
(<http://root.cern.ch/root/>)
 - GANGA implementado en Python
(<http://ganga.web.cern.ch/ganga/>)

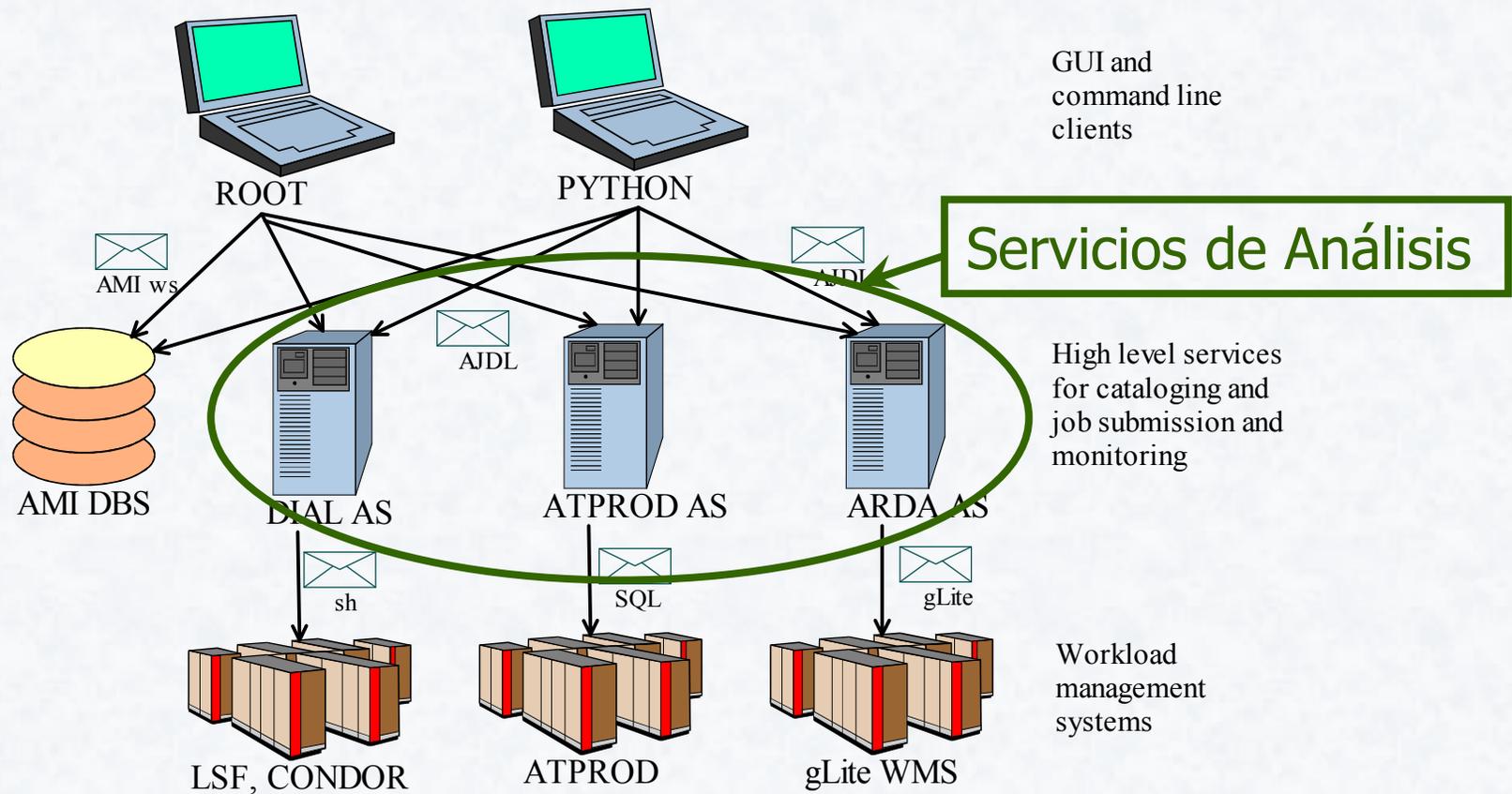
Servicio de catálogos



Servicio de catálogos

- ✓ Proporciona catálogos de los objetos de C++ de AJDL
- ✓ Permite a los usuarios:
 - Asignar metadata (características) a dichos objetos
 - Registrar la procedencia de los datasets
 - Monitorizar el progreso de los trabajos
- ✓ Tres tipos de catálogos:
 - Repositorio → para todo tipo de objetos de AJDL
 - ❑ Almacena descripciones de objetos en formato XML
 - Catálogo de selección → para datasets
 - ❑ Asocia metadata con los objetos
 - Catálogo de réplicas → para datasets
 - ❑ Asocia un nombre lógico con una colección de réplicas
- ✓ Actualmente se usa una base de datos de tipo MySQL como servicio de catálogos

Servicios de análisis



Servicios de Análisis

✓ Objetivo

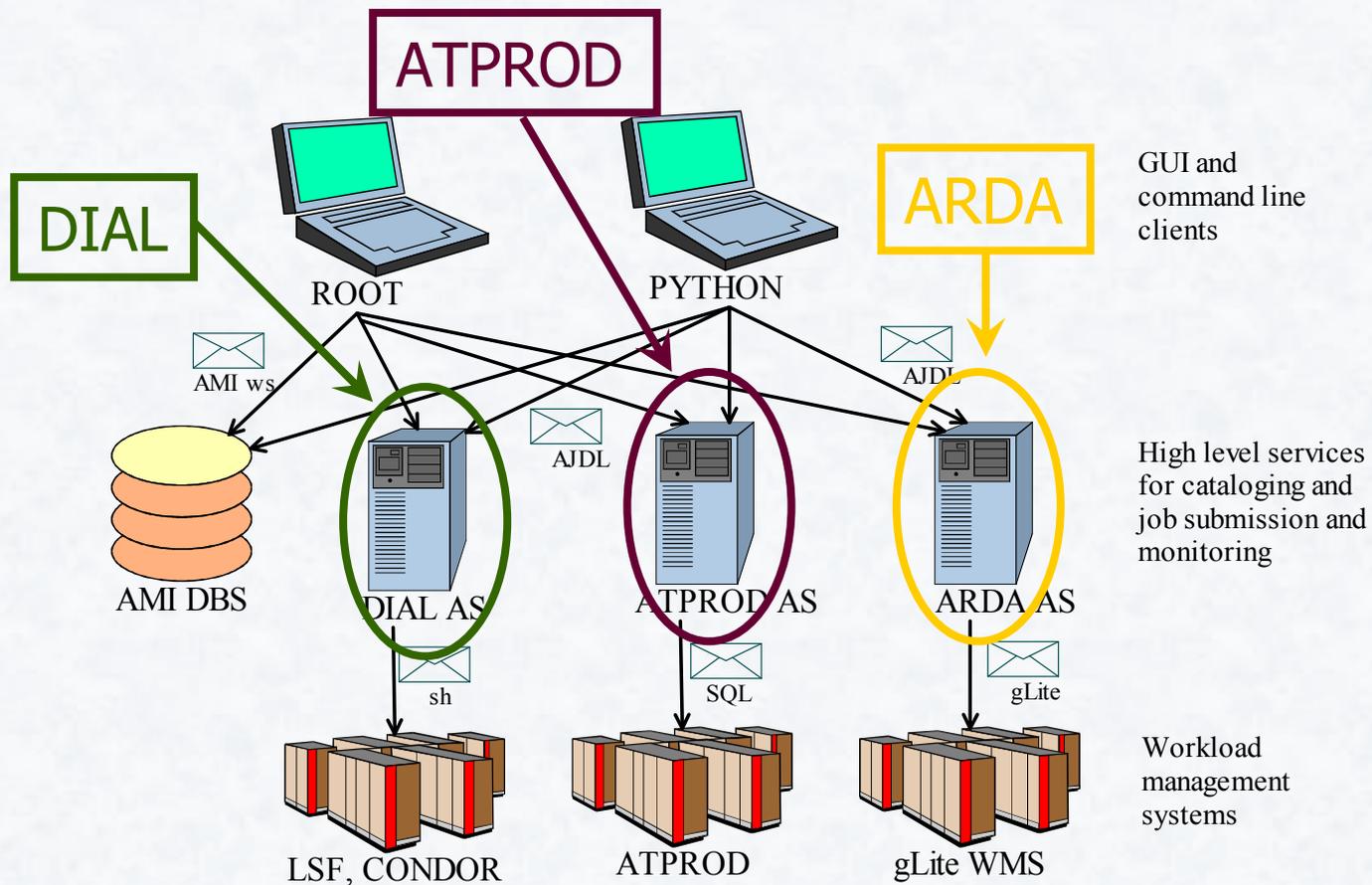
- Proporcionar al usuario los mecanismos necesarios para facilitar la gestión de los trabajos
 - métodos para la instalación de transformaciones, el envío de los trabajos y su monitorización.

✓ Procedimiento

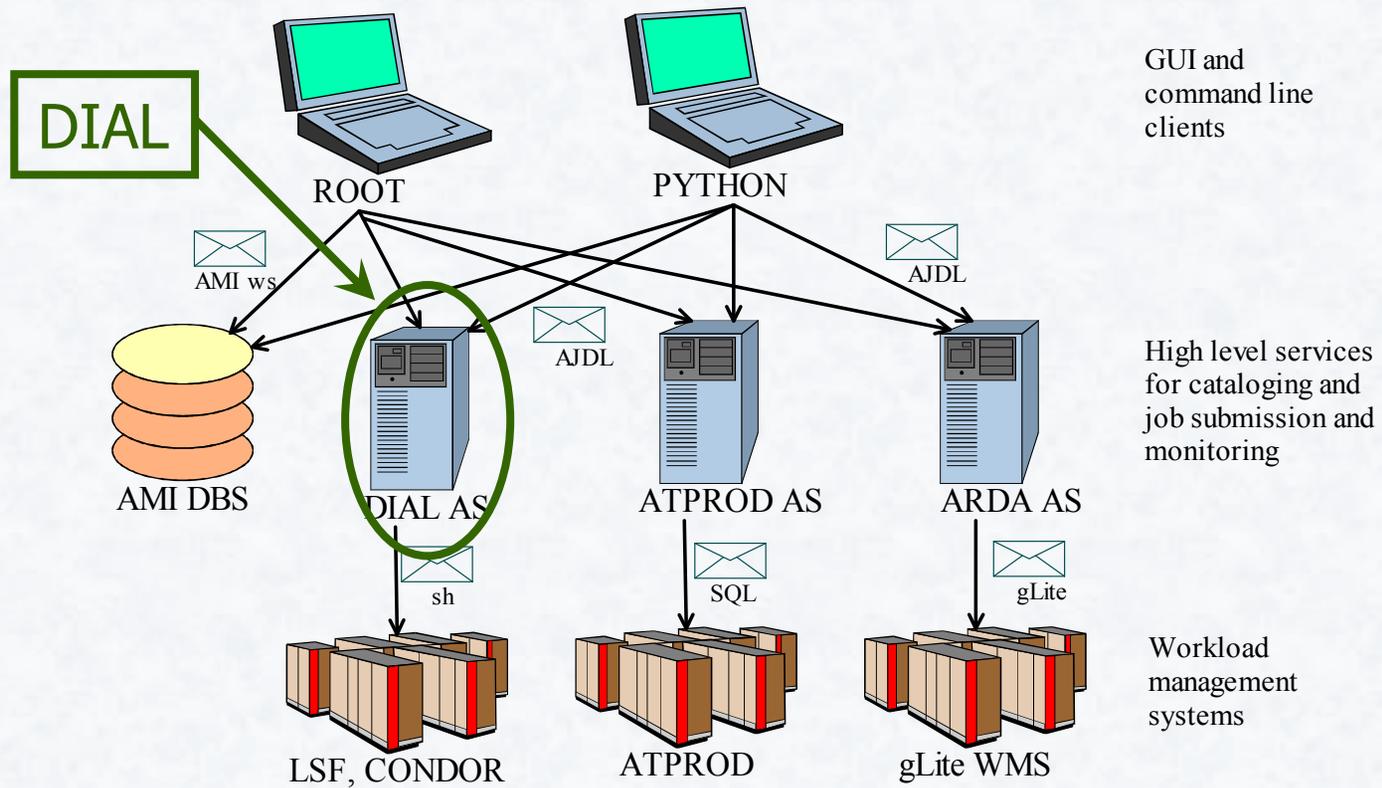
- Recibe una petición de trabajo desde AJDL
- Crea el correspondiente trabajo asignándole un identificador
- Envía el trabajo para su ejecución al sistema de procesamiento

✓ Actualmente ADA posee tres servicios de análisis:

Servicios de análisis



DIAL



DIAL

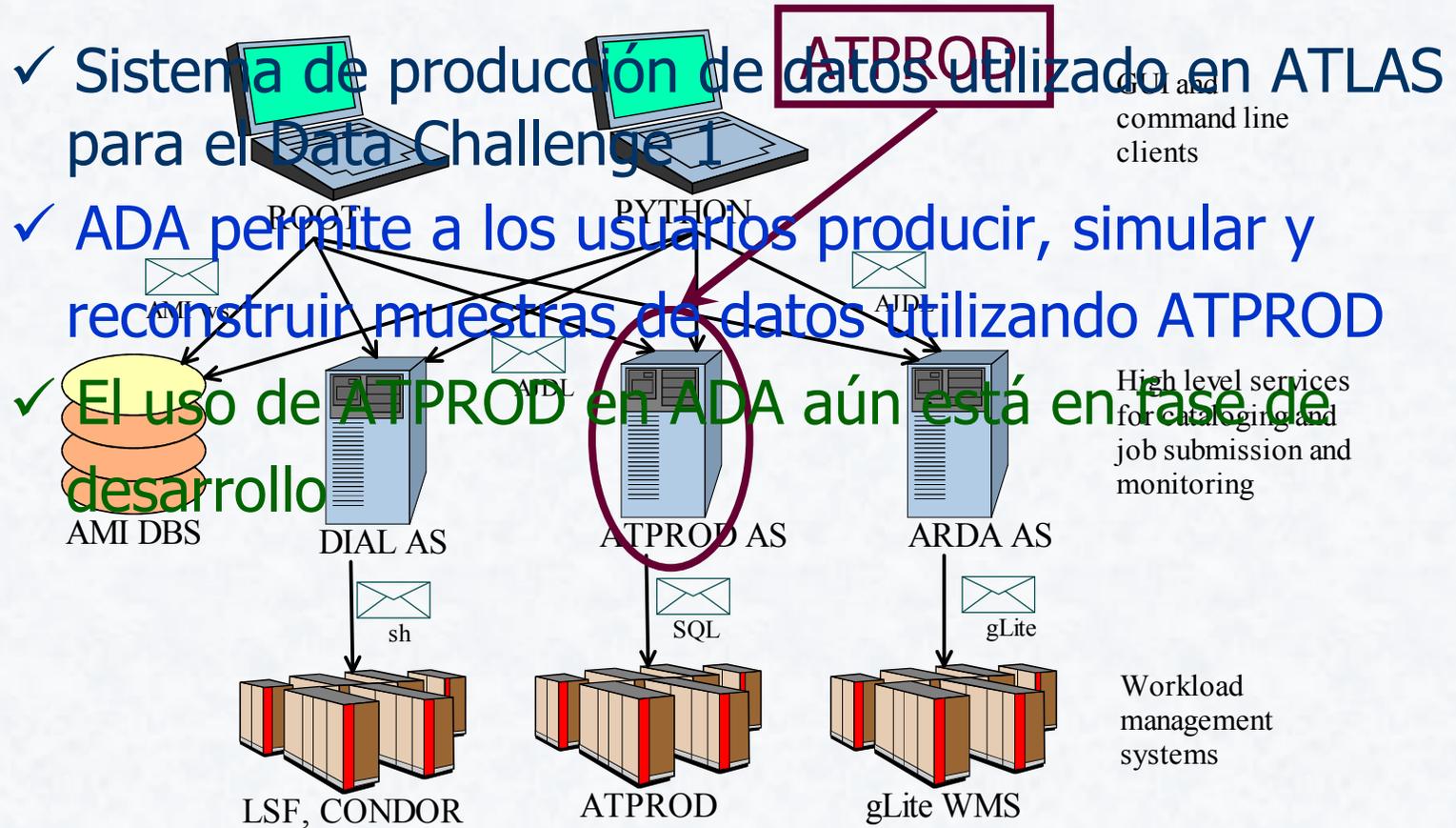
- ✓ Contiene la librería de objetos de C++ definidos bajo el nombre de AJDL
- ✓ Demuestra la viabilidad del análisis distribuido e interactivo de grandes cantidades de datos, y el no-interactivo de los mismos
 - análisis interactivo → colas rápidas (bajo LSF)
 - análisis no-interactivo → colas lentas (bajo Condor)
- ✓ Su componente principal es el Scheduler que: --->
 - recibe una petición del usuario (transformación + dataset)
 - divide el dataset en sub-datasets
 - por cada sub-dataset crea un sub-trabajo
 - los envía al sistema de procesamiento para su ejecución
 - concatena los diferentes resultados dando el resultado final

ATPROD

✓ Sistema de producción de datos utilizado en ATLAS para el Data Challenge 1

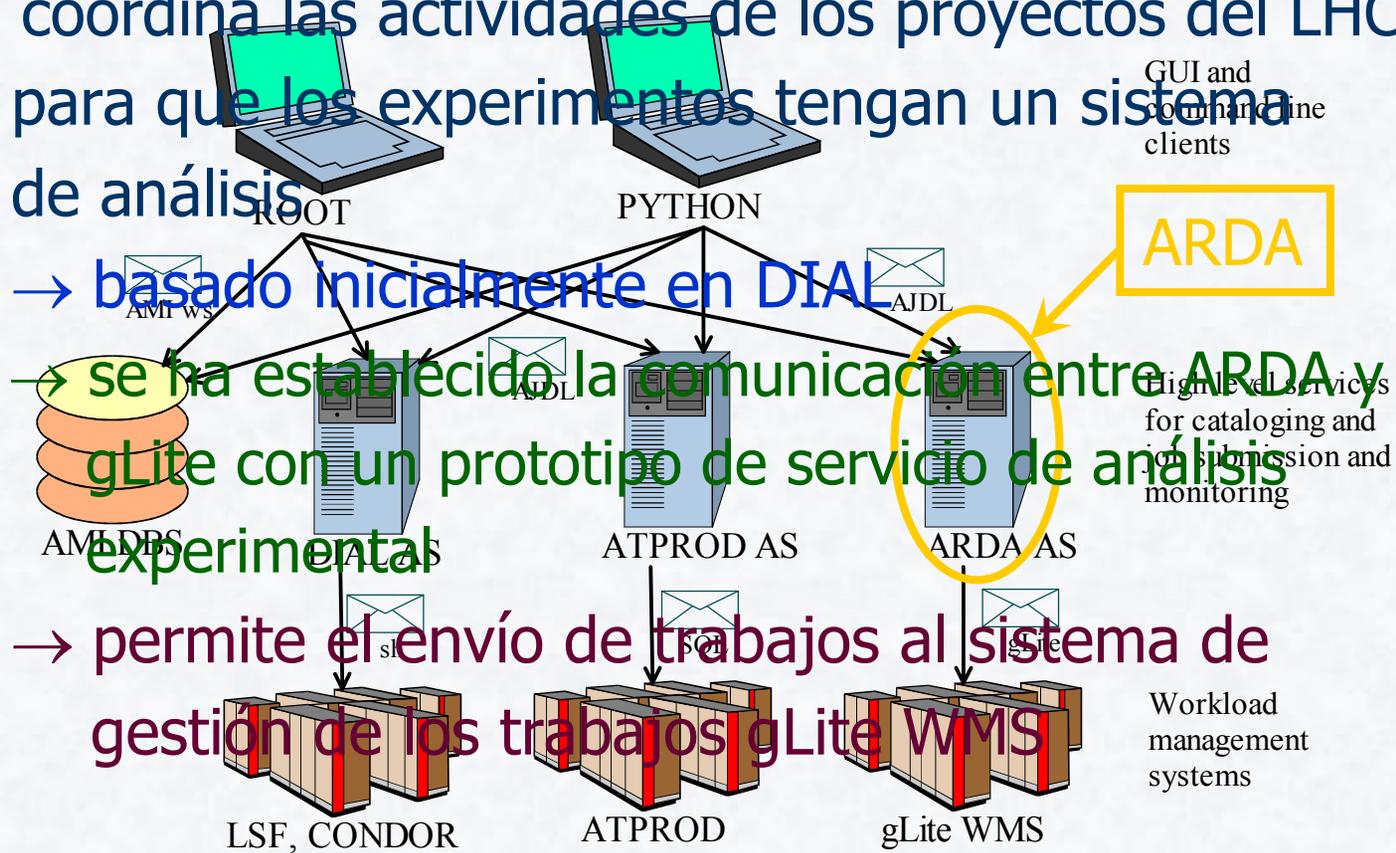
✓ ADA permite a los usuarios producir, simular y reconstruir muestras de datos utilizando ATPROD

✓ El uso de ATPROD en ADA aún está en fase de desarrollo



ARDA

- ✓ coordina las actividades de los proyectos del LHC para que los experimentos tengan un sistema de análisis



- basado inicialmente en DIAL
- se ha establecido la comunicación entre ARDA y gLite con un prototipo de servicio de análisis experimental

- permite el envío de trabajos al sistema de gestión de los trabajos gLite WMS

Conclusiones

- ✓ El Sistema de Análisis Distribuido de ATLAS facilitará a los físicos de ATLAS el procesamiento y el análisis de los datos que se producirán en dicho experimento
- ✓ Admite un amplio rango de aplicaciones pero se centra en aquellas que permiten el análisis, la generación, la simulación y la reconstrucción de grandes cantidades de datos
- ✓ El análisis y el procesamiento de los datos se realiza de manera completamente distribuida
- ✓ En fase de desarrollo