# Long Term Data Preservation for CDF at INFN-CNAF

S.Amerio[1], L.Chiarelli[2], L.dell'Agnello[3], D.De Girolamo[3], D.Gregori[3], M.Pezzi[3],A.Prosperini[3], P.Ricci[3], F.Rosso[3], D.Salomoni[3], S

[1] University of Padova,[2] GARR, [3] INFN-CNAF

## INTRODUCTION

Interest in the long term preservation of scientific data and their availability to general public is growing. Data collected in High Energy Physics (HEP) experiments are the result of a significant human and financial effort. The preservation of HEP data beyond the lifetime of the experiment is of crucial importance to ensure the long term completion and extension of scientific programs, to allow cross collaboration analysis, analyzing data from several experiment at once, to perform new analysis with new theoretical models and techniques and for education, training and outreach.

HEP data preservation poses many technical and organizational challenges: data preservation implies migration to new storage media when available, adjusting data access methods if needed; data analysis capabilities must be preserved ensuring the experiment legacy software runs on new platforms, or on old ones with no security issues; validation systems have to be set up to regularly check data access and analysis framework; all the information needed to access and analyze data has to be properly organized and archived.

## CDF Computing Model

CDF ended its data taking in 2011, after collecting 10 fb$^{-1}$ of data.

**Data**: 5 PB, stored on T10K technology tapes at Fermilab in three formats: Raw, Reconstructed, User level's ntuples. Metadata stored in Oracle databases. Data handling based on SAM (Sequential data Access via Metadata) and dCache.

**Software**: Code (C, C++, Python) in frozen releases in CVS repositories. Latest version for SL5; a SL6 legacy release ready by the end of 2013. External dependencies: GEANT3, CERNLIB, Neurobayes, Root, Oracle.

**Job submission**: CDF *Central Analysis Farm code (CAF)* provides the users with a uniform interface to resources on different Grid sites.
Three portals based on glideinWMS to access computing resources at Fermilab, OSG, Tier1 @ CNAF and other LCG sites.
Authentication:
• Kerberos + FNAL KCA to obtain a X.509 certificate
• CNAF VOMS to setup a valid proxy that can be used to submit to the grid.

## CDF Data Preservation Project at INFN-CNAF

Goal: preserve a complete copy of CDF data and MC samples at CNAF + services (access, data analysis capabilities)

The project is divided into two main areas:

**Copy of the data**
The copy will be splitted in two years
• end 2013 - early 2014 → All data and MC user level ntuples (2.1 PB)
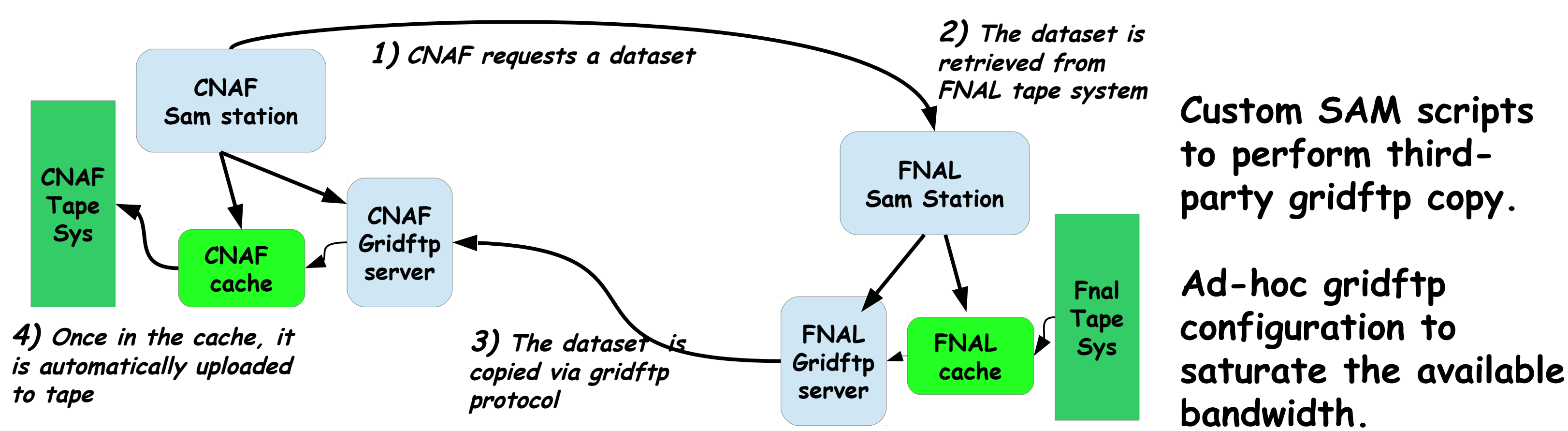• mid 2014 → All raw data (1.9 PB) + Databases

**Long term future analysis framework**
Based as much as possible on the current resources already available at CNAF (Job submission portal, SAM station)
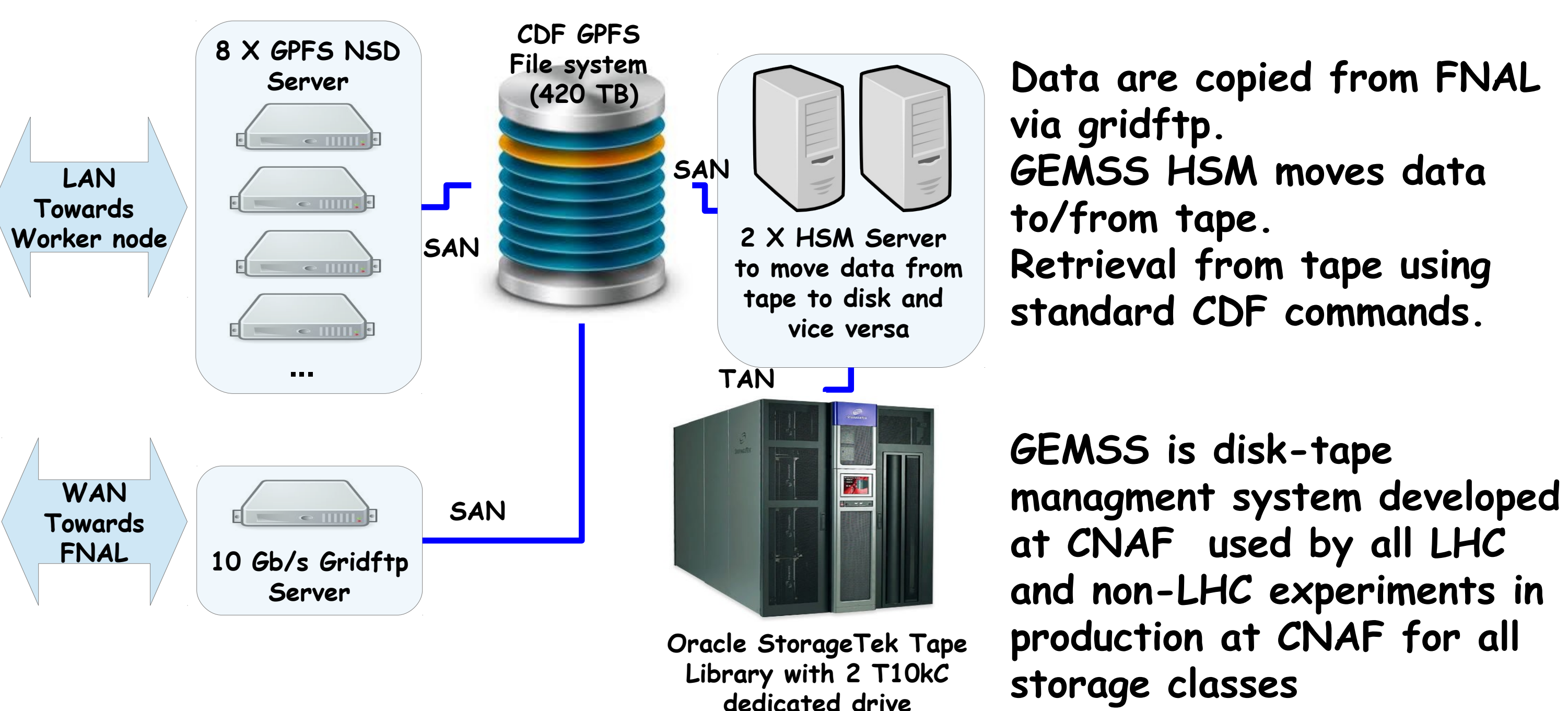To run CDF legacy code we plan to use a dynamic virtual infrastructure through the INFN-developed WNoDeS framework.
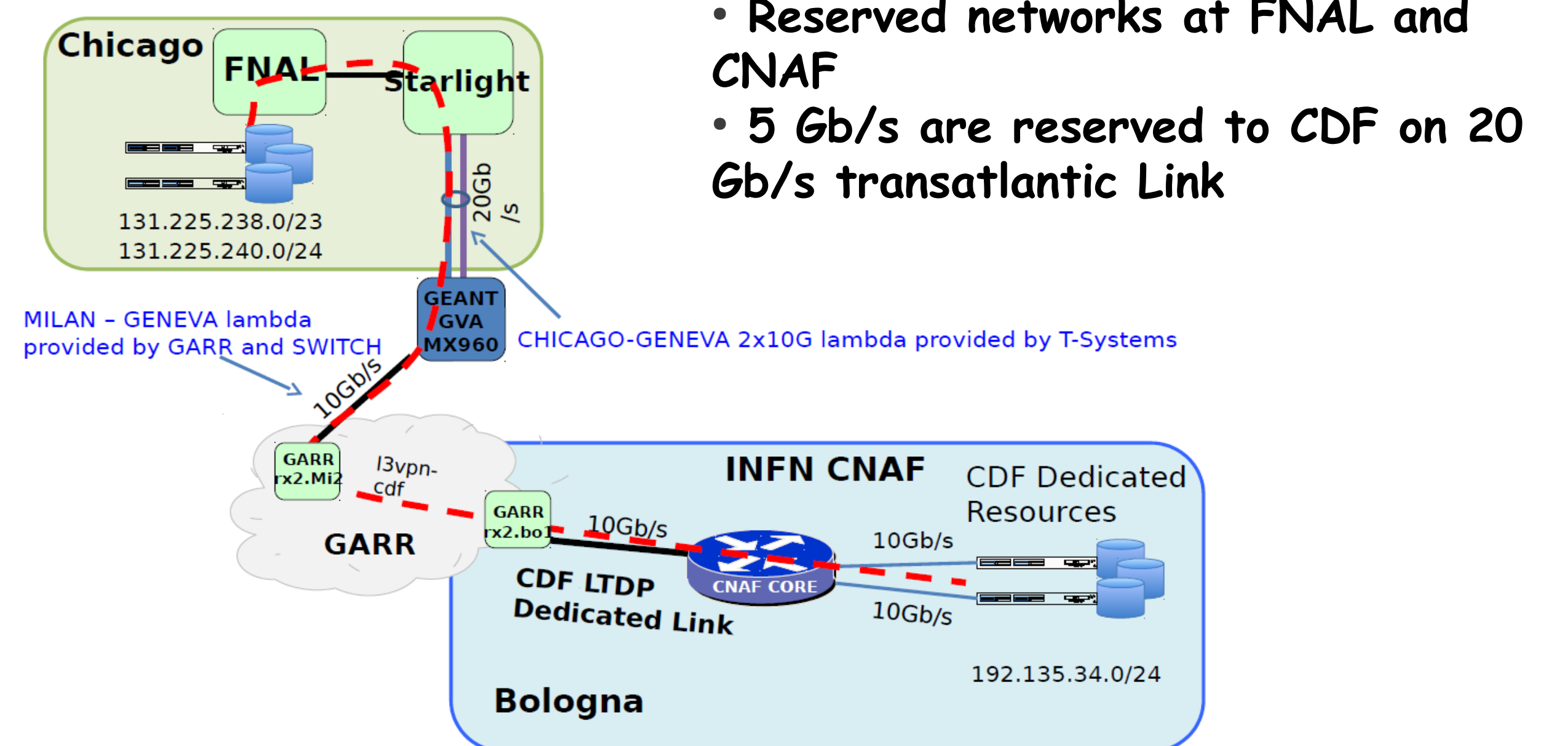
## CDF data at CNAF
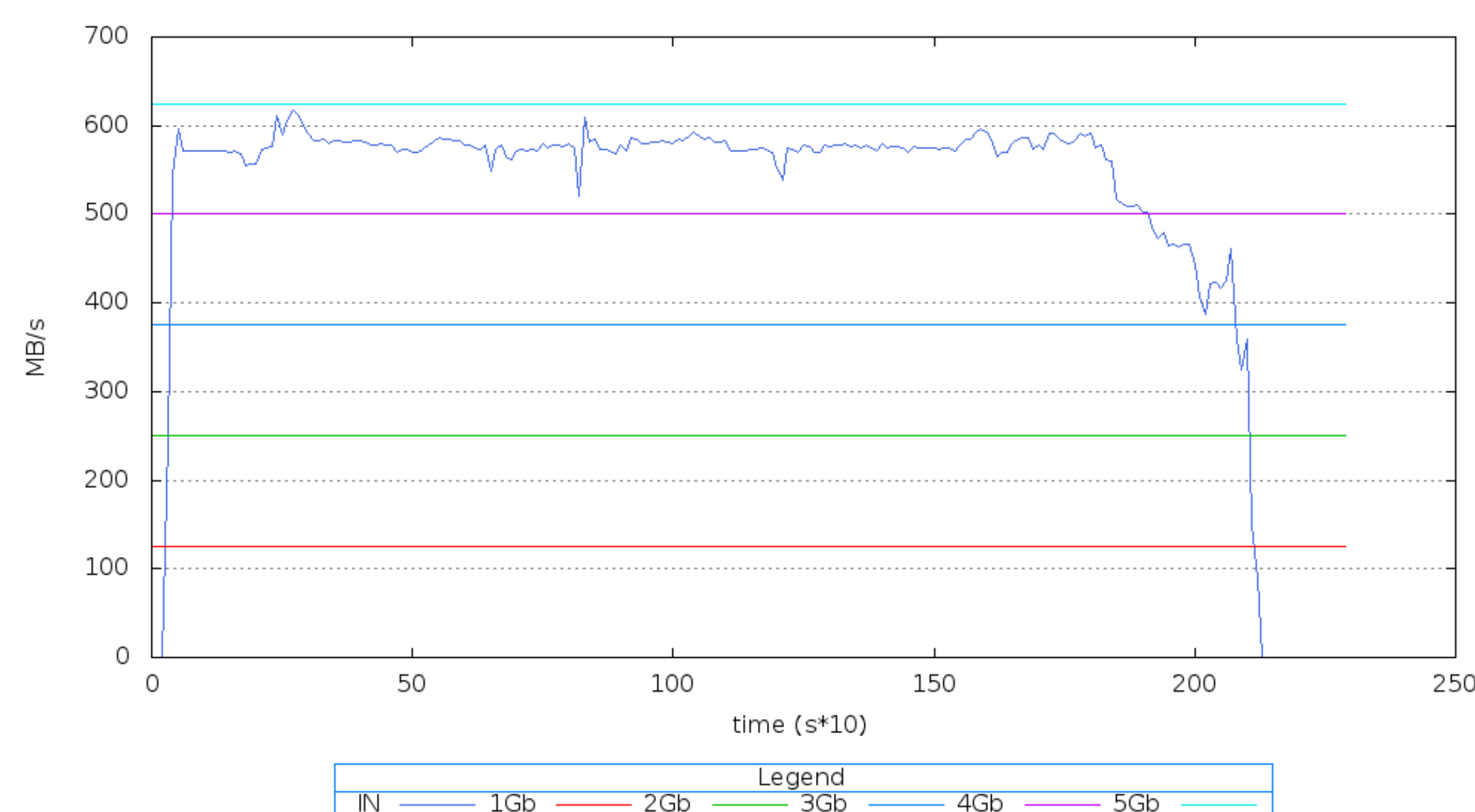
### CDF FNAL-CNAF Data Copy Mechanism



Custom SAM scripts to perform third-party gridftp copy.

Ad-hoc gridftp configuration to saturate the available bandwidth.

### Network Layout for FNAL-CNAF CDF Data Copy



• Reserved networks at FNAL and CNAF
• 5 Gb/s are reserved to CDF on 20 Gb/s transatlantic Link

### Storage Layout for FNAL-CNAF CDF Data Copy



Data are copied from FNAL via gridftp.
GEMSS HSM moves data to/from tape.
Retrieval from tape using standard CDF commands.

GEMSS is disk-tape managment system developed at CNAF used by all LHC and non-LHC experiments in production at CNAF for all storage classes

### Data transfer rate FNAL-CNAF



Data transfer rate stable over time

With ~ 50-80 parallel copy processes we exploit at the best the available bandwidth..

## CDF data analysis in the long term future

To analyze CDF data stored at CNAF:
• SAM station to retrieve data from tape: new version of SAM code in preparation at FNAL.
• Disk buffer to stage data retrieved from tape
• User area
• CDF code volume, accessible via AFS
• Access to Databases
• Computing resources to run CDF legacy code

**WNoDeS → Worker Nodes on Demand Service**

In production at several Italian centers, including the INFN Tier-1 since November 2009.
Dynamic virtual networks, *new* feature under development:
dynamic instantiation of private VLANs and address assignement for VM isolation.

In the long term future, CDF services and analysis computing resources can be instantiated on demand on pre-packaged VMs in a controlled environment.