

CDF Computing Politics and Technicalities

Stefano Belforte
INFN - Trieste

Talk content

- It is a hard (2) days' night
 - Will show summary of each slide only
 - Show the real slide only on demand
- Much more to know and to say than can fit here, see e.g.:
 - <http://www.pd.infn.it/CSN1/riunioni/28-01-2002>
 - <http://www.pd.infn.it/CSN1/riunioni/24-06-2002>
 - <http://www.pd.infn.it/CSN1/riunioni/16-09-2002>
 - http://www.ts.infn.it/~belforte/offline/index_offline.html

Talk content

- Overview and critical history
- Status and perspective
- CDF vs. CDF-Italy
- CDF vs GRID
- Impact on INFN
 - CSN1
 - CNAF
- Much more to know and to say than can fit here, see e.g.:
 - <http://www.pd.infn.it/CSN1/riunioni/28-01-2002>
 - <http://www.pd.infn.it/CSN1/riunioni/24-06-2002>
 - <http://www.pd.infn.it/CSN1/riunioni/16-09-2002>
 - http://www.ts.infn.it/~belforte/offline/index_offline.html

History

- CDF offline upgrade was planned as Run1 extrapolation

History

- 3 guidelines for CDF computing upgrade Run1→Run2
 - All new code, all new hardware. A big thing !
- 1. Build on Run1 success
 - Data was analyzed
 - No major drawback emerged
- 2. Smooth introduction of C++
 - Allow wrapped Fortran and "banks" to survive for a while
- 3. Fix most acute problems
 - Data access
 - ☞ hand mounted tapes
 - ☞ scripts with lists of file names
 - Bookkeeping
 - ☞ reproducibility of past results
 - ☞ offline version + calibration constants

Italy vs CDF

- Italy not a part in CDF offline upgrade
 - Main Italian contribution to CDF:
 - ☞ Detectors - Trigger
- We have no faults
- But will claim some merits

Italy vs CDF

- Italy not a part in CDF offline upgrade
 - Main Italian contribution to CDF:
 - ☞ Detectors - Trigger
- Code contribution (Padova, Rome) to high level analysis, but...
 - No charge in management, design, infrastructure, major code development, DataBase, farming...
- Only "management" role: SB as internal reviewer and chair of "computing forum"
- SB, head of computing for Italy, had SVT as first priority until 1 year ago and spent 5 years full time on it
- We do not like the way offline has gone and wish we had more impact, but the only way to be heard is by working

CDF Sociology

- We could not/can not impact the basics of CDF computing
- Only way to direct the cart is by pulling yourself
- So far we did not pay either

CDF Sociology

- CDF is a loosely coupled organization
- Very few and limited institutional responsibility in software, besides FNAL
- Most work is done on a voluntary basis
 - People need visibility (to get a job e.g.) or must write their Ph.D. thesis
 - Hard to "direct". Quality uncertain
- FNAL committed to provide basic software infrastructure
 - Small group, no strong leadership
 - Bottom line: if it runs, is enough
- Only way to direct the cart is by pulling yourself
- Especially since we did not put any money in the offline upgrade itself

Computing Hardware Responsibilities

- No formal allocation of responsibilities as yet
- By unwritten agreement FNAL had always provided everything
- So far lot of people added their (private) own anyhow
 - Including us
 - This works
- FNAL considering of asking CDF Institution for sharing of computing cost because it was suggested from outside
- None in CDF is advocating sharing of FNAL costs, nor shared usage of offsite resources

Computing Hardware Responsibilities

- No formal allocation of responsibilities as yet
- By unwritten agreement FNAL had always provided
 - Data storage and Event Reconstruction ("Production")
 - Analysis facilities available for free to everybody
 - ☞ FNAL project → external input difficult
- It sort of worked
 - Several people did analysis using only FNAL resources
 - Several analysis (including Italians) relied on significant non-Fnal hw at home or at the lab that was:
 - ☞ unshared, poorly coupled to data repository
- FNAL considering of asking CDF Institution for sharing of computing cost because it was suggested from outside
- None in CDF is advocating sharing of FNAL costs, nor shared usage of offsite resources

How the original plan evolved

- New code runs very slow (OO ?)
- Biggest hit:
 - analysis hardware needs, underestimated by x100
 - only place where a radical change was required
- Most other stuff worked though

How the original plan evolved

- OO proved less friendly than previous stuff
 - Slow learning curve
 - Code harder to read
 - Banks→objs: Documentation from poor to none
 - ☞ the hacker's motto "use the force, read the source"
 - CPU needs for simple tasks increased x10
- The "Run1 model" for analysis hardware architecture broke
 - Plan had to change from few big SMP's to 1000 PC's
- But other parts of offline upgrade did very well:
 - event reconstruction (production farm) OK
 - Data Handling OK (heavily revised, but little extra \$)
 - bookkeeping OK
 - ☞ Oracle "has its price" though: \$\$ and complexity

CDF computing

- Data hierarchy (~500TB/year) :
 - DAQ 30Hz average
 - Event 10⁹/year
 - DataSet ~ 1TB files from one trigger path
- Data format: Root I/O
- Data Handling built around data organization:
files with metadata
 - Transparent access to tapes
 - Data access by dataset name

CDF computing

- Data hierarchy (~500TB/year) :
 - DAQ 75Hz peak, 30Hz average
 - Event $10^9/\text{year}$ raw/dst/ana = 250/400/100 KB
 - File 1GB smallest unit Data Catalog knows of
 - DataSet N files from (selection of) one trigger path
1TB typical
- Data format
 - Root I/O (random access possible), not Root objects
 - Root multibranching (in same file!) planned for 2003
- Data File Catalog: Oracle
- Data Handling built around data organization (files with metadata) and user's needs:
 - Transparent access to tapes
 - Data access by dataset name

Production at CDF is not a problem

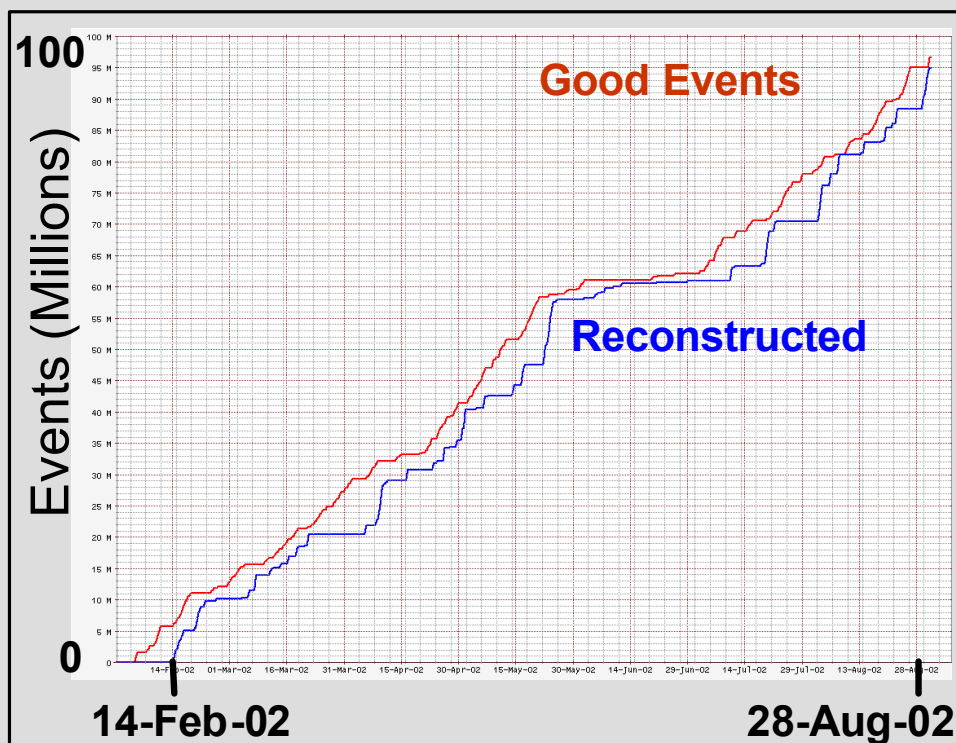
- Full reconstruction: **~4sec/event** on present CPU
- 150 nodes farm is enough
- No hardware crisis here
- Can't imagine a reason to "distribute reconstruction around"

Production at CDF is not a problem

- Full event reconstruction: **~4sec/event** on 1GHz P3 CPU
 - a bit less than **2K SpecInt2000*sec / event**
- 4x30Hz = need 120 CPU to keep up with average DAQ rate
 - Add reprocess, some MC, handle rate ~ peak (75Hz):
 - ☞ have 170 duals now (about half "old" <1GHz)
 - ☞ equivalent to 300 x 1GHz CPU
 - Still a small farm, similar to CDF Level 3
 - 3 persons run it + 1FTE for hw support + 1 shifter
- No hardware crisis here, in spite of OO
 - Run1 number about 5x lower (~ 700MIPS*sec)
 - ☞ "then" estimate for Run2 (higher lum): 1200 MIPS
 - ☞ difficult to judge the remaining x3 (new COT/SVX)
 - Have more or less (x4 ?) same size farm now and then
 - Did not take advantage of Moore's law, but hw is cheaper
- Can't imagine a reason to "distribute reconstruction around"

Farms Reconstruction Progress

- Farms keeping up with data and have lots of reprocessing capacity
 - 95 million events reconstructed at least once in Feb-Aug-02
 - Processing ~6 million event/day as we speak (70Hz vs 75DAQ)



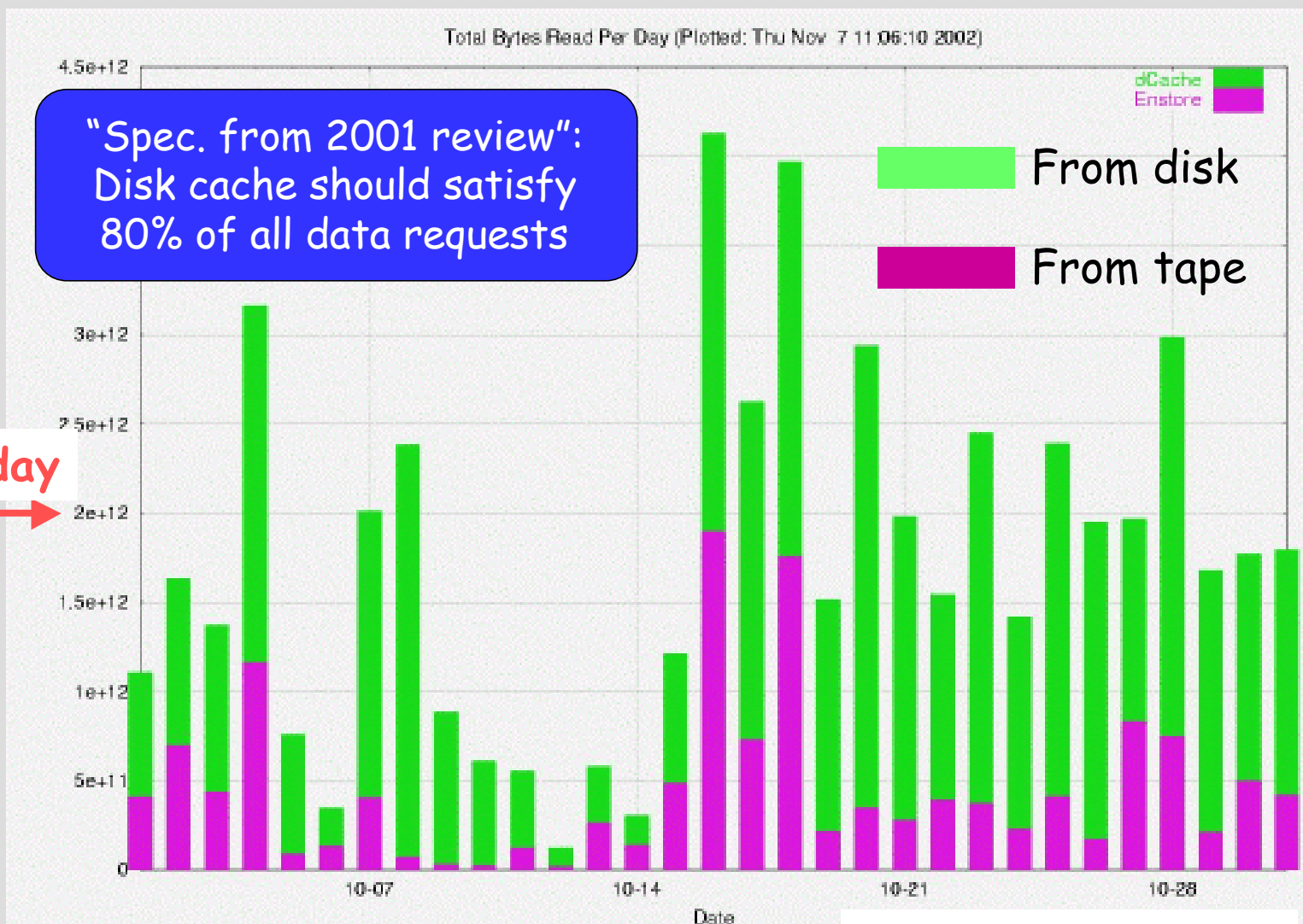
CDF Data Handling

- DH works well
- Transparent access to tapes now (even from Italy)
- Hierarchical distributed disk cache coming

CDF Data Handling

- **Disk Inventory Manager**
 - CDF home made product, based on DataFileCatalog
 - **Intelligent disk cache** (when a dataset is requested, hands to user files already on disk first)
 - **Works** (after switching from Sony AIT to STK)
 - But **only on one machine**, hard to use offsite or in farm
- **Enstore+dCache: local (Fnal) disk cache** in front of networked tape handling (STK)
 - It works (even from Italy)
 - Access by file name. DIM style **optimization missing**
- **SAM**: another oracle DB, **supercedes** present **DataFileCatalog**
 - **Distributed hierarchical disk cache on WAN**, recovers DIM intelligent approach to using cached files first
 - Talks to Enstore, will integrate dCache
 - It is **evolving into GRID** product

Tape to Disk to CPU



Days in October 2002

Stefano Belforte - INFN Trieste
CDF computing

Analysis Hardware

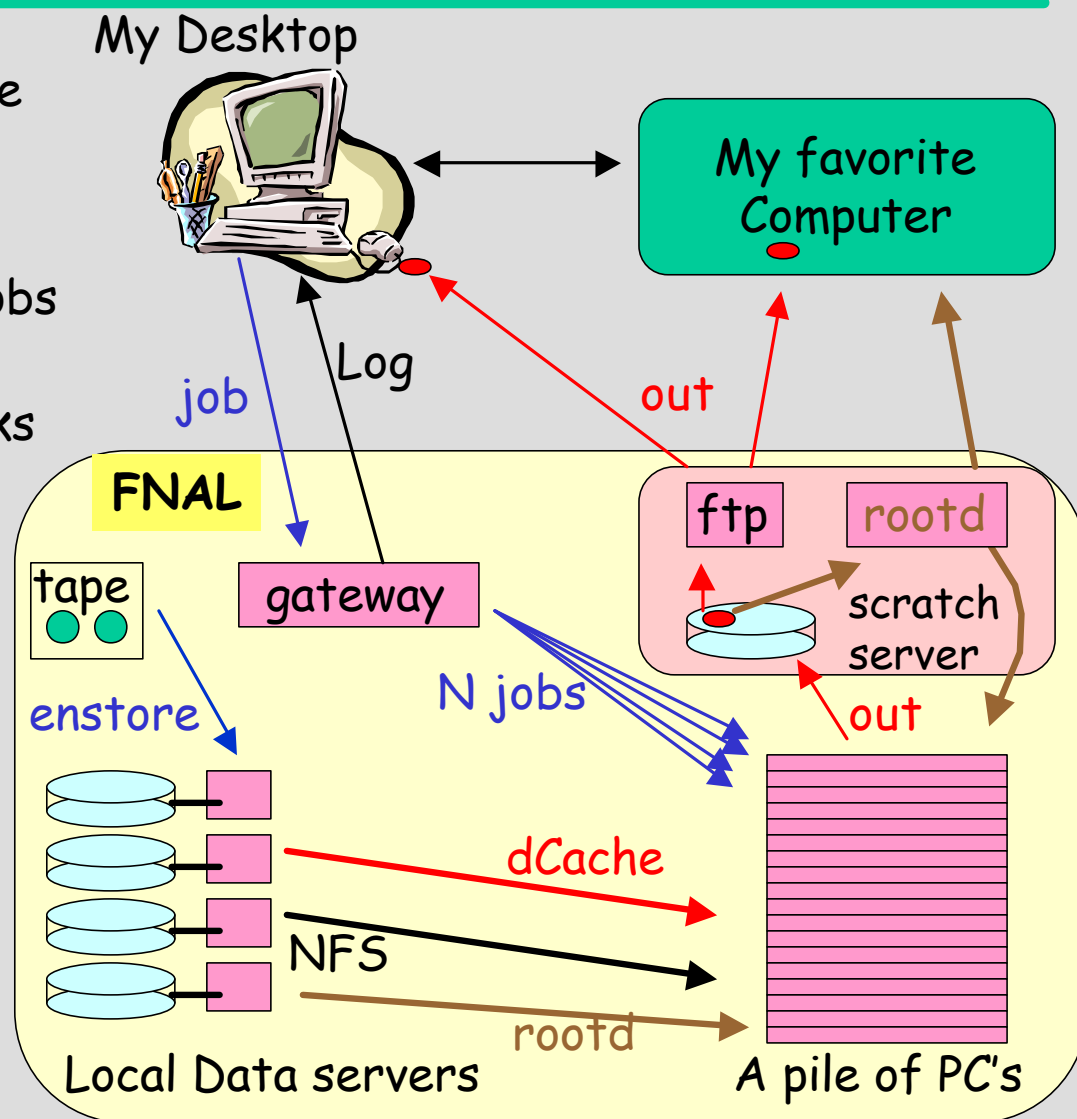
- 2001 review
 - Need a thousand disk drives and a thousand PC's
"The new CAF"
 - Will give 100x the previously planned capacity (5 years later) at the same cost
- New CAF built because MIT put 3 full time persons on it
- Italy put a lot of work too

Analysis Hardware

- This is where OO hit us hard
 - User's code runs at ~1MB/sec, has to go through 1TB
- Code development was not driven by speed optimization
- The Run1 model was very expensive hardware, differently from production farm: SMP and FC RAIDs in a SAN
- When it came to need 100x the Run1 CPU... it did not scale
- Extensive review of model and needs in fall 2001 lead to:
 - Need to navigate 1TB in few hours (100 parallel jobs)
 - Need disk cache O(100)TB (a thousand disks)
 - Need analysis farm O(1000) GHz (a thousand PC's)
 - Not an easy conclusion. The old plan had reasons: I/O.
- Problem is: new big farm is a hardware nightmare and no demonstrated solution yet. Especially disk access.
- Was only done because MIT group put 3 FT persons on it
- Italy also put a lot of work on batch and monitor

CDF Central Analysis Farm

- Compile/link/debug everywhere
- Submit from everywhere
- Execute @ FNAL
 - Submission of N parallel jobs with single command
 - Access data from CAF disks
 - Access tape data via transparent cache
- Get job output everywhere
- Store small output on local scratch area for later analysis
- Access to scratch area from everywhere
- **IT WORKS NOW**
- Remote cloning in progress



Fermilab budget for analysis

- 2M\$ spent on analysis hardware before 2002
(out of 10M\$ total)
- 3M\$ in 2002-5 for the CAF requested (will get less)
 - 2003: cpux10 diskx10
 - 2005: cpux100 diskx50

Fermilab budget for analysis

- Overall estimate in 1997: 10M\$ - all spent by 2002
- Request for next years: 2M\$/year (most likely will get >1.5)
- Analysis hardware, about 2M\$ spent before 2002
 - ~1.5M\$ in one big SMP: 128*300MHz SGI
 - ☞ CPU ~40GHz i.e. equivalent to 40 "2001 PC's" (50k\$)
 - Plus ~40TB on FC SCSI RAID at about 2x IDE cost
 - Move to CPU farm + RAID on IDE disk servers
 - ☞ 2.5K\$ per dual P3 2GHz
 - ☞ 12K\$ per 2TB disk server } 2002 prices
 - ☞ requested about 0.8M\$/year in 2002-5 for new CAF
 - ☞ will have 10x the SGI CPU next year (for the same \$)
 - ☞ will have 100x the SGI power by end of 2005
 - ☞ will have 500TB data disk by end of 2005

Strategy for the Italian CDF group

- Analyze Ntuple where you are
- Make Ntuple at FNAL and copy to Italy as needed
- Referees helped in making correct decisions
- We stayed into our 1999 budget estimate, in spite of "slow code" crisis
- Now we have our private share of CAF at FNAL
- decentralizedCAF in progress: our CAF in Italy ?

Strategy for the Italian CDF group: 1

- **Interactive work** (analysis of Ntuple):
 - **Desktops/Mini-farms** in Italy
 - **Desktops at FNAL**
 - ☞ mostly small (~2.5KEu)
 - ☞ few "powerful" (~3.5KEu) (full time residents e.g.)
 - ☞ will explore farm solution with CAF team
- **Plan evolved significantly in last 3 years** (but as we promised at the very beginning, we kept overall cost for Run2a ~3GLit) :
 - **SMP → farm transition**
 - ☞ When needs escalated, we changed model, not requests
 - **Very productive interaction with our referees**
 - ☞ they helped us to slow down, think hard, make the best decision at the proper time, stick to the budget

Strategy for the Italian CDF group: 2

- Started as: have hw at FNAL to make ntuple, copy those home
 - "Computers at FNAL" OKed by CSN1 (with 1MEu cap)
 - Referees did not like the SMP's, nor FC, nor SCSI
 - SMP's and SAN killed anyhow by CPU crisis
 - Italians worked to understand the problem and needs, to indicate an alternative solution and make it work
 - Now the Central Analysis Farm (CAF)
 - ☞ Italians buy a piece of it and get priority of usage
 - ☞ Working well as we speak
 - ☞ Minimum effort, maximum gain:
completely managed by CDF, optimal access to data
- Moving to adding offsite decentralized farms (dCAF)
 - To make it easy to use offsite computers
 - INFN also pushing us to have computers in Italy

Data Handling from Italy

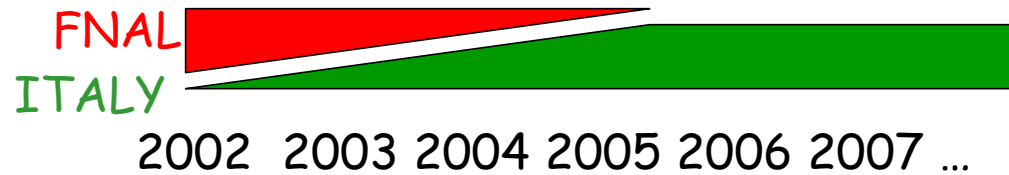
- CDF problem is access to data
- Working on Italy on **static copies** makes **little sense**
- Need local transparent disk cache: SAM
- Working in Italy **less efficient because of distance from tape repository**
- There is no technical reason to put hardware in Italy rather than at FNAL

Data Handling from Italy

- The CDF problem is access to data
- Run2a: 1PB of primary data, 200TB of analysis data (PADs)
 - Growing x8 while moving from Run2a to Run2b
- Working on Italy on **small static copies** makes **little sense**
- **Replicating all data** in Italy makes **even less** sense
- Working mechanisms for local caching and a good network are the way. WAN not a problem anymore. Caching is.
 - SAM may be the solution, R&D in progress
- **Working in Italy will never be as efficient as FNAL, due to distance from tape repository**

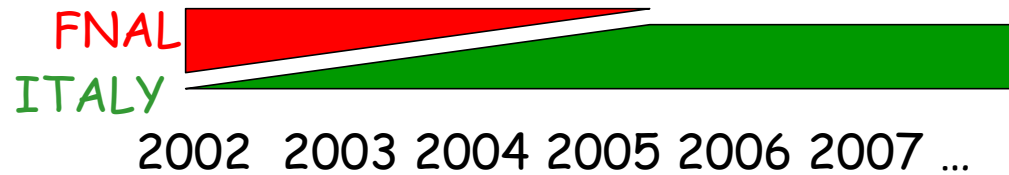
- There is no technical reason to put hardware in Italy rather than at FNAL

Timeline



- 2002: batch at FNAL, interactive in Italy
- 2003: batch at FNAL: interactive in Italy
 - Start batch analysis in Italy
- 2004: try “all in Italy”, but do not rely on it
- 2005: if all goes well, leave FNAL

Timeline



- **2002: batch at FNAL, interactive in Italy**
 - Test presence at CNAF waiting for infrastructure
- **2003: batch at FNAL: interactive in Italy**
 - Demonstrate CNAF on a few simple realistic analysis
 - First significant hardware purchase at CNAF for CDF
 - Test CNAF as provider of services
 - Test usage of GRID tools for transparent access
- **2004: try "all in Italy", but do not rely on it**
 - Demonstrate CNAF on large analysis
 - Replicate at CNAF processing capability from FNAL
 - Test CNAF as provider of smooth 24x7 operations
- **2005: if all goes well, leave FNAL**
 - Keep expanding CNAF x2 every year

Caveats for a CDF analysis farm at CNAF learn from the CAF

- An analysis farm is much more than a pile of PC's
- If we clone CAF it will work
 - But CDF will be different from other EDG based stuff
 - ➔ it will require dedicated support (besides physicists)
- If we go another way (EDG ?)
 - ➔ it will very likely require more people for less results

Caveats for a CDF analysis farm at CNAF learn from the CAF

- An analysis farm is much more than a pile of PC's
- If we clone CAF, from batch system (FBSNG) to authentication (Kerberos5 in FNAL.GOV realm) to job submission (CafGui) to local data caching (SAM)
 - we know it will work
 - can build on FNAL experience and trade knowledge
 - it will require dedicated support (besides physicists)
 - startup (2~3 months) likely 2 FTE
 - after startup: 1FTE sw + "whatever it takes" for hw
- If we go another way (EDG ?) we do not know what we face, we may not be able to reuse CDF tools, we have no support for farm performance tuning, remote data access, I/O, batch optimization, handling of user's data, job submission from FNAL (no EDG UI there)
 - it will very likely require more people for less results

Moving CAF to CNAF: the proposal

- Will try and see, decision to leave FNAL will have to be based on proof of existence of valid alternative here
- Our needs from CNAF may be very different from an LHC-Tier1 role, we will not feel bad if left out, but very upset if pulled in and then set aside w/o support or resources

Moving CAF to CNAF: the proposal

- Start with limited, but significant hardware
 - 2003 at CNAF $\approx \frac{1}{2}$ of private share of CAF in 2002
 - Present estimate: "after summer 2003"
 - 5TB of disk and 20 dual processor agreed at last CSN1
 - ☞ based on estimate of work on one analysis stream
 - How to avoid duplication of resources in 2004/5 ?
- Explore effectiveness of work environment
 - Don't give up on CAF features
 - Look for added value
 - Will need help (manpower)
- Will try and see, decision to leave FNAL will have to be based on proof of existence of valid alternative here
- Our needs from CNAF may be very different from an LHC-Tier1 role, we will not feel bad if left out, but very upset if pulled in and then set aside w/o support or resources

The CDF Italy Plan for Analysis (june 24)

- Estimates of CPU and disk requirements for analysis by the italian CDF group presented to CSN1 on June 24 2002
- Satisfying needs for up to 18fb^{-1} (full CDF Run2)
- Including 40% contingency
 - 2002-2004 at FNAL: 1M Euro
 - 2005-2008 at CNAF: 2M Euro
 - 2002-2004 interactive in Italy : 0.5M Euro
 - MC and interactive at FNAL → contingency

The CDF Italy Plan for Analysis (june 24)

year	Luminosity		ANALYSIS FARM			contingency 40% (Keuro)	Requested per year (Keuro)
	Planned (Church)	Target (adjusted)	disk (TB)	CPU (duals)	cost/y (Keuro)		
2001	commissioning		0.6	0			43
2002	0.3	1.0	20	80	336	0	336
2003	1.2	2.0	40	140	266	106	372
2004	2.5	3.5	70	200	285	114	399
TOTAL cost Analysis Farm at FNAL + 40% conting. for Run2a expanded to 3.5 fb-1							1150
2005	4.1	6.0	110	280	331	132	463
2006	7.6	9.5	180	350	298	119	417
2007	11.3	13.5	250	430	331	132	463
2008	15	18.0	330	500	288	115	403
TOTAL cost for Analysis Farm at CNAF + 40% coting. for Run2b (15 fb-1)							1746
TOTAL BUDGET CENTRALIZED COMPUTING FOR ANALYSIS 2001-2008							2896

- ❖ Only analysis farm. No MC. But 40% contingency next years.
- ❖ Will cover up to 3.5 fb-1 with money indicated last year for 2
- ❖ Future farm at CNAF may cost up to a factor 2 more while dealing with 5 times the data in 4 instead of 3 years.
- ❖ Overlap of resources during FNAL/CNAF transition not included

What do others do

- Join the CAF
 - Investment per person similar to ours
- Private resources at FNAL
 - Sometimes very large
- Stuff at home
 - Several US groups have large resource at universities
 - Outside US: very small groups, very little hw, except
 - ☞ UK: 1.8M pounds (~3M Euro) already
 - ☞ Canada: 220 duals cluster already

What do others do

- **Join the CAF**
 - 1 fileserver/4persons + 4duals/person typical
 - INFN 2002: 7 fileservers, 60 duals (~30 persons)
- **Private resources at FNAL**
 - MIT's 30 nodes MOSIX
 - Small clusters ~1TB ~5duals
 - High end dual (~4.5K\$) on each desk
- **Stuff at home**
 - Several US groups... have large resource at Universities shared with other groups, not counted in CDF budget, managed to do a lot of remote analysis in the past
 - Outside US: very small groups, usually little hw, but...
 - ☞ UK: 1.8M pounds (~3MEuro) already in hand
 - ☞ Canada: cluster of 220 duals already there

SPARES

From now on:
spare slides

I have (almost) not talked about Monte Carlo

- Nor I will
- It is supposed to be a low impact job
- Never a big need in the past
 - Everybody struggled just to make what was needed for his/her analysis
- Nobody pulling the cart of a large common initiative, it may just end up in the same way as the past
- Biggest problem is good bookkeeping (as usual) from scripts to random numbers
- Anyhow:
 - Reconstruction farm has spare cycles for 60MEv/year
 - Canada just offered Toronto's 448 CPU's (300MEv/year)
 - Mostly waiting for physicists
- If it comes to "do your own" Italian needs will likely be covered within the contingency indicated in the plan

CDF-FNAL budget requests for next years

CDF Plan and Budget for Computing in Run 2

Version 3
May 16, 2002

Edited by
Robert M. Harris
Fermilab Computing Division

Contributions from
William Badgett, Stefano Belforte, Phil Demar, Richard Jetton, Kevin McFarland, Don Petravick, David Tang, Jeff Tseng, Steve Wolbers and Frank Würthwein

1. Start with needs as function of integrated luminosity

Requirement	Offset	Slope
Batch CPU (GHz)	66 (GHz)	1100 (GHz/fb ⁻¹)
Static Disk (TB)	32 (TB)	125 (TB/fb ⁻¹)
Read Cache (TB)	12 (TB)	35 (TB/fb ⁻¹)
Write Cache (TB)	5 (TB)	10 (TB/fb ⁻¹)
Disk I/O (GB/s)	368 (MB/s)	1100 (MB/s/fb ⁻¹)

2. Compute overall HW requirements

FY	02	03	04	05	06	07	08
Batch CPU (THz)	0.51	1.5	2.9	4.7	8.6	12	17
Farm CPU (THz)	0.37	0.70	0.76	1.3	2.1	2.5	2.8
Static Disk (TB)	82	180	340	540	980	1400	1900
Read Cache (TB)	26	54	100	160	280	410	540
Write Cache (TB)	9	17	30	46	81	120	160
Disk I/O (GB/s)	0.81	1.7	3.1	4.9	8.7	13	17
Archive Volume (PB)	0.3	0.7	1.1	1.7	2.9	4.1	5.3
Archive I/O (MB/s)	65	190	130	480	350	470	590

Table 3: The integrated computing needs by the end of each fiscal year as a function of fiscal year.

3. CDF budget requests to CDF/FNAL/DOE

FY	Batch CPU (\$M)	Inter. CPU (\$M)	Farm CPU (\$M)	DB (\$M)	Tape Robot (\$M)	CAF Disk (\$M)	Cache Disk (\$M)	Net-work (\$M)	Legacy Sys. (\$M)	Total (\$M)
2002	0.59	0.07	0.22	0.02	0.77	0.47	0.16	0.25	0.69	3.24
spent	(0.19)	(0.07)	(0.11)	(0.00)	(0.25)	(0.12)	(0.04)	(0.12)	(0.69)	(1.59)
2003	0.48	0.15	0.22	0.15	0.35	0.35	0.11	0.25	-	2.06
2004	0.48	0.2	0.13	0.10	0.35	0.35	0.11	0.25	-	1.97
2005	0.60	0.2	0.19	0.10	0.35	0.35	0.13	0.25	-	2.17

Table 18: CDF computing equipment purchasing plan. The fiscal year, batch CPU for the CAF, interactive CPU and its local disk, production farm CPU, databases, tape robot & tape drives, network attached disk for the CAF, read and write cache disk, networking, legacy CPU and disk systems, and total procurements.

What do others do at FNAL ?

- Joining the CAF
 - INFN : 7 filesystems 60 duals
 - Pittsburgh : 8 duals
 - Carnegie-Mellon : 8 duals
 - KEK-Japan : 2 filesystems 38 duals
 - Korea : 0.5 filesystem (+ 2 later)
 - Spain : 1 filesystem
 - Canada : 1 filesystem
 - Switzerland : 1 filesystem
 - UK : 15 filesystems (most for common use)
 - More US (8 universities) : 10 filesystems 4 duals
- Having their private stuff at FNAL
 - MIT : 30 nodes Mosix cluster
 - Many US: small clusters ~2TB ~5 duals (Duke, Rochester)
 - Glasgow : 10 duals

What other US groups do at home

- Hard to say, a couple examples:
- TexasTech (6 people): ~40CPU beowulf cluster
 - cloning SAM and CAF now
- Rutgers (4 people): ~6duals
 - Cloning SAM and CAF now
- Chicago (10people): probably 20 PC's, few TBs
- Carneige-Mellon, LBL, Pennsylvania, Argonne, Illinois,... have large resource at universities, managed to do a lot on remote analysis in the past, how much they can use of those common resources is probably undefined until afterwards. Those computer are not payed by DOE money and never appear in any CDF "chart".

CDF computing outside US (approx)

	2002		2003		Notes
	TB	duals	TB	duals	
Spain	-	-	10	50	Shared with CMS, plan for EDG tools No plan for shared access
Germany	3 + 1	20 + 10	20 + 20	50 + 40	Tier1 (shared with LHC) + Tier3 (CDF) No plan for shared access Testing SAM on Tier3
UK (4 sites)	24	16	80	64	Maybe 5x the CPU if 8-way → duals No EDG, Kerberos for user access, SAM for data. <i>maybe open</i>
Korea	1	20	7	40	Want to clone CAF by end of 2002 Kerberos for user access, <i>open to all</i> Start w/o SAM
Canada	1	8	28	224	No GRID tools <i>Run official CDF MC and copy to FNAL</i>
Italy	1	5	5	20	No plan for shared access Exploring SAM on single node