

Dalla presa dati alla pubblicazione @ *BaBar*

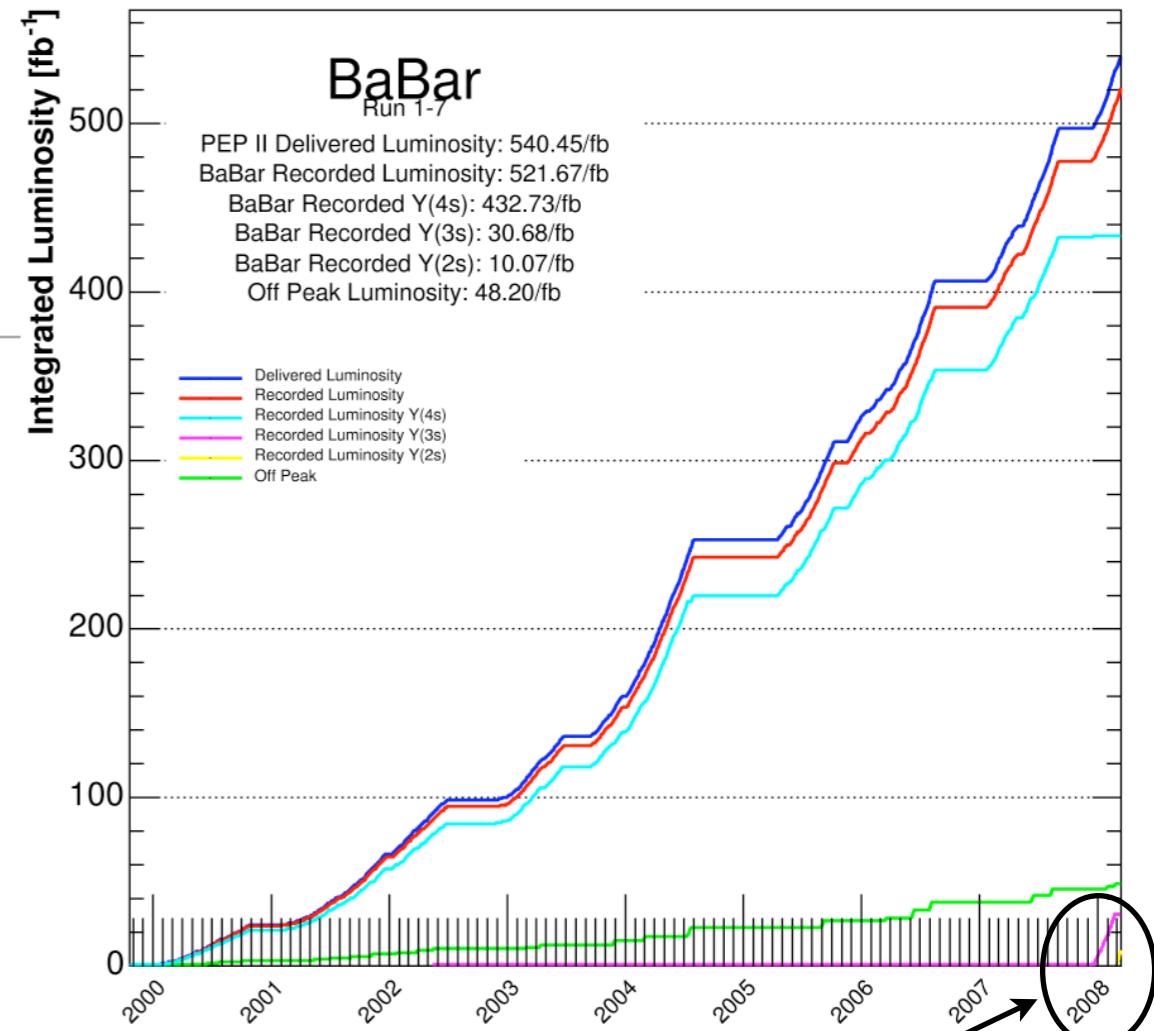
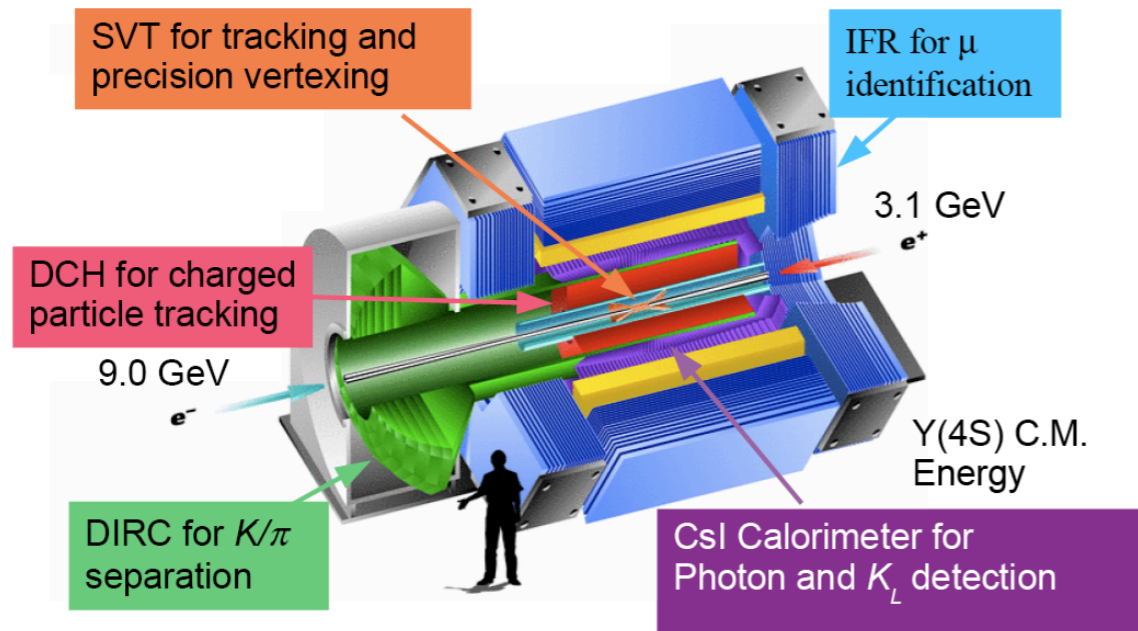
Alfio Lazzaro
Milano Atlas Meeting
18 Marzo 2008

ATLAS



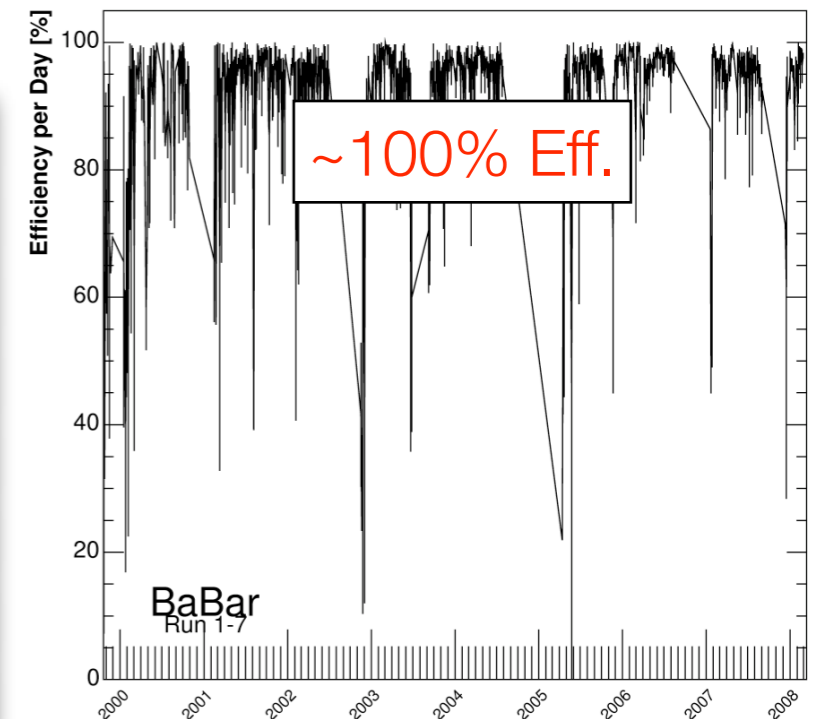
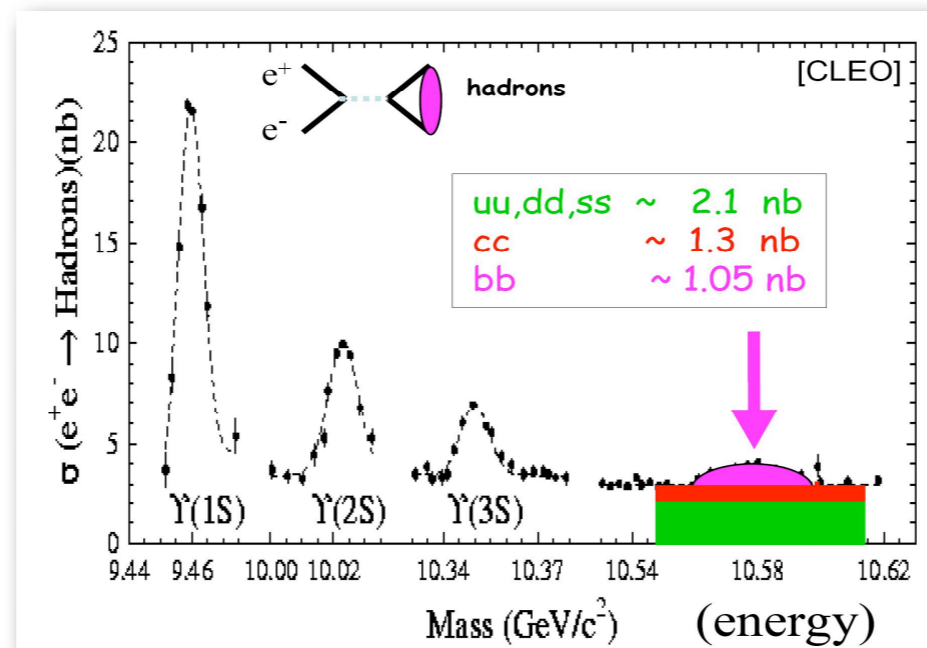
BaBar @ SLAC

- Operativo dal 1999
- Energia nel CM: ~ 10.58 GeV

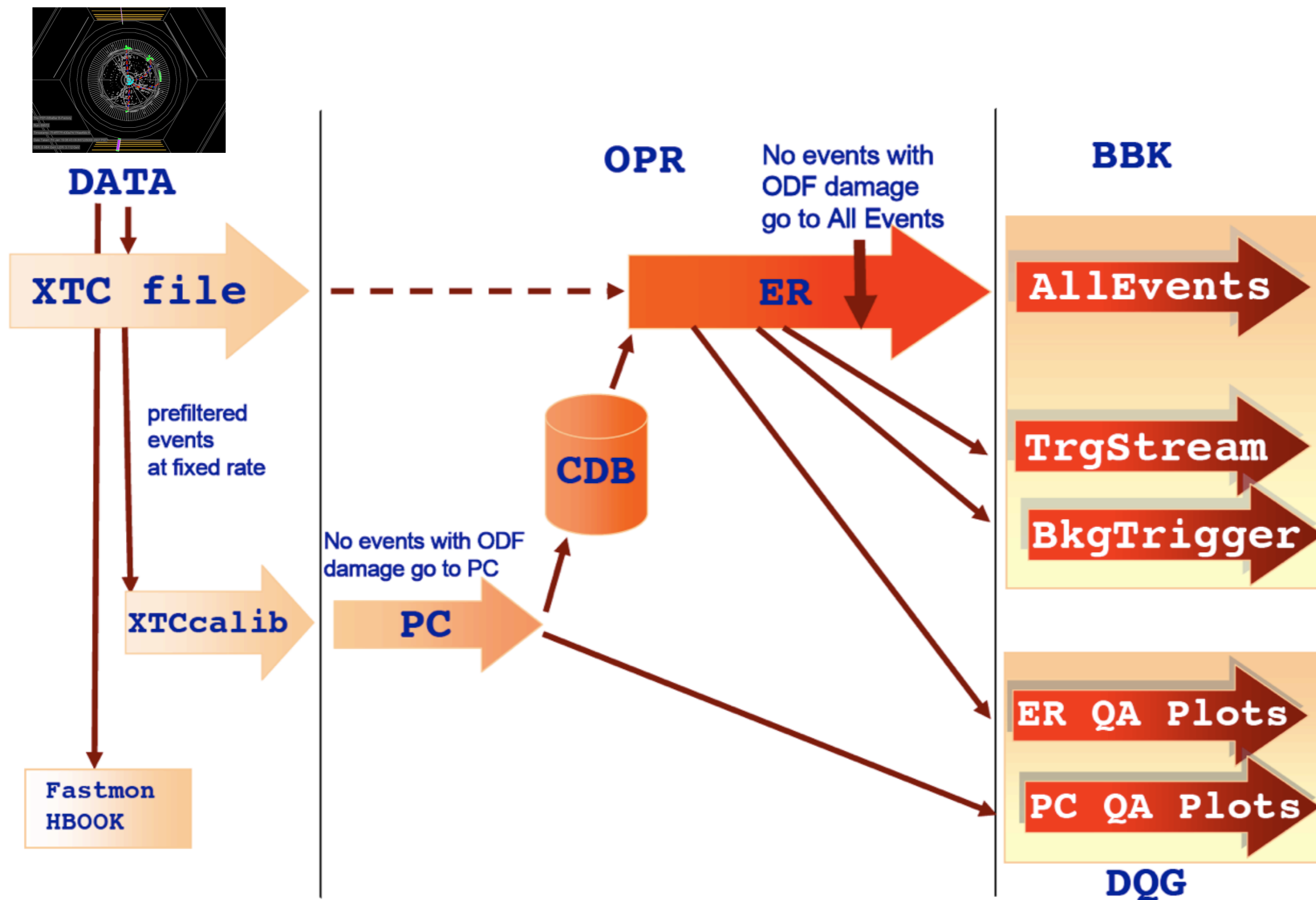


$\Upsilon(3S)$ & $\Upsilon(2S)$

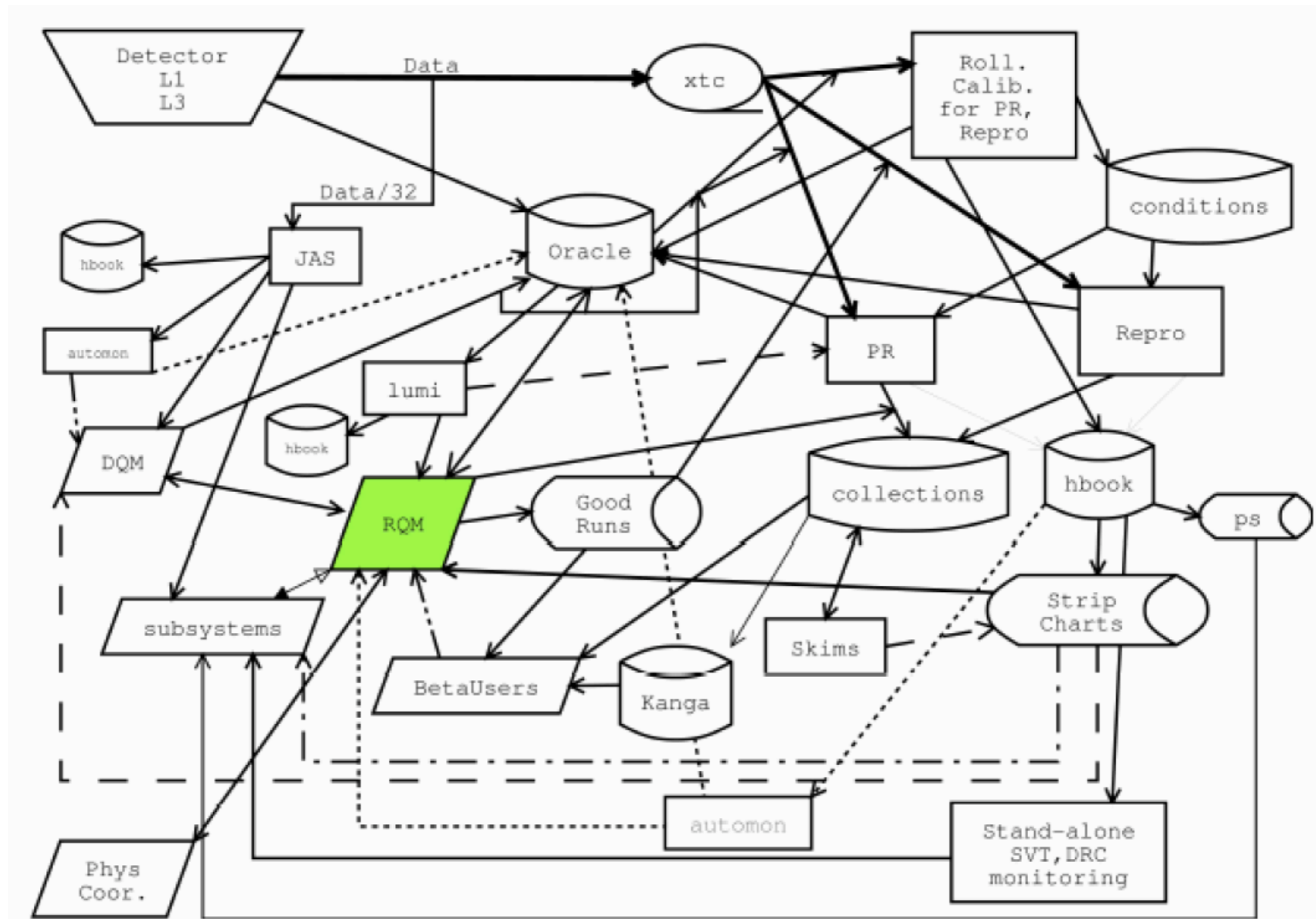
- Ambiente molto pulito:



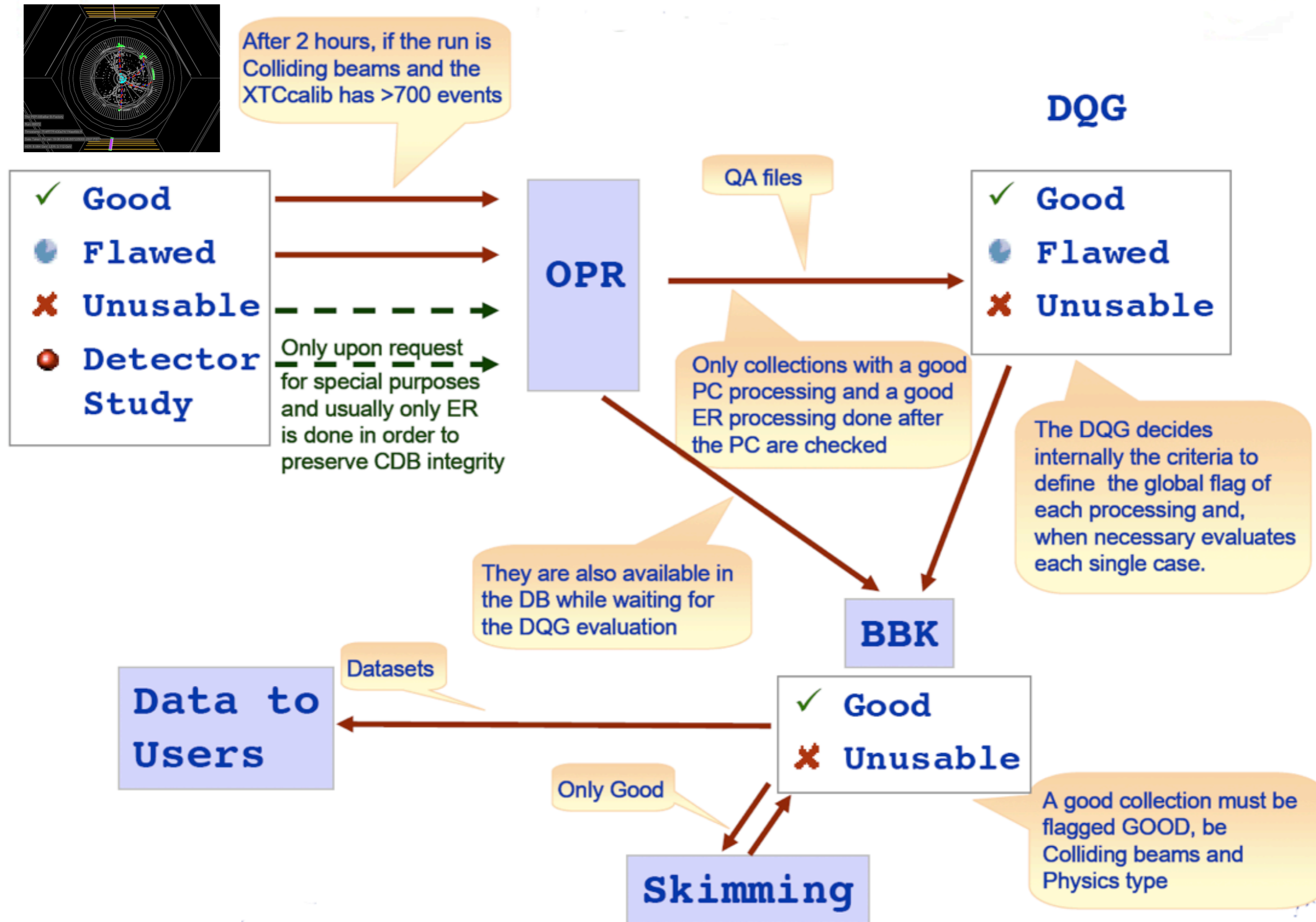
I dati: dall'acquisizione al riprocessamento (mondo ideale)



I dati: dall'acquisizione al riprocessamento (odissea reale)



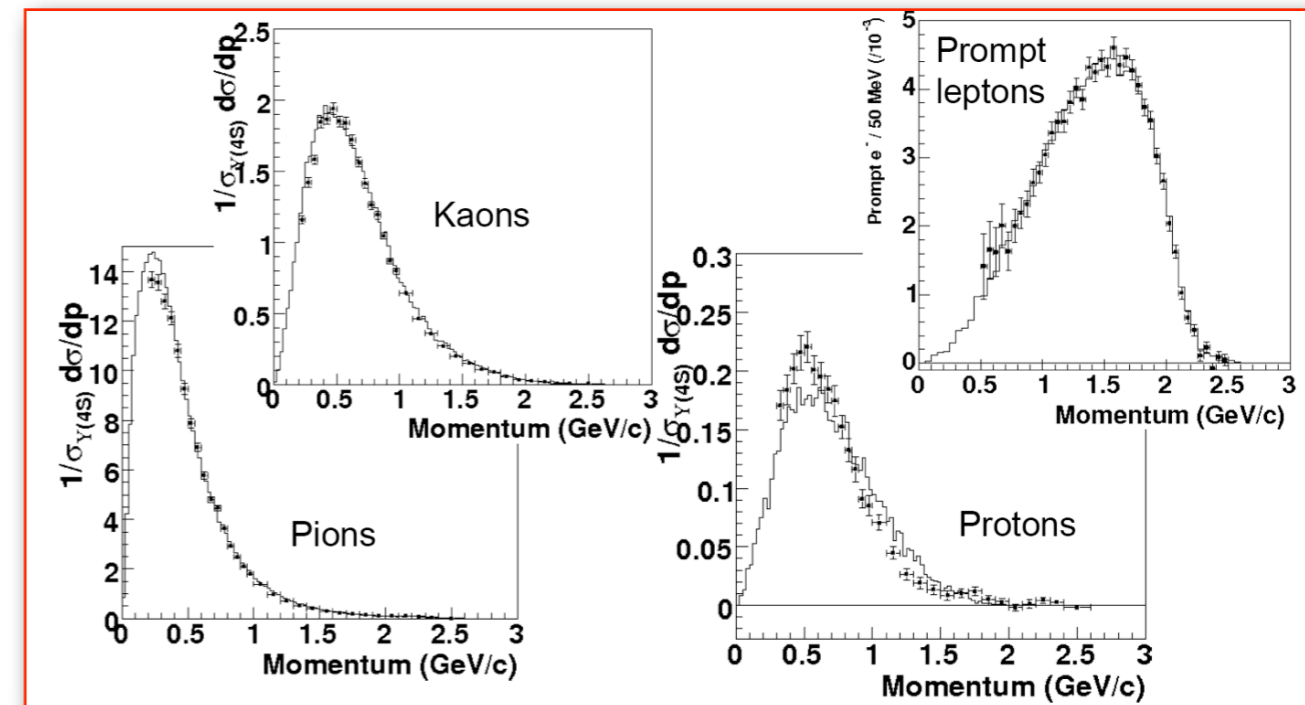
I flags



La simulazione

- Diversi generatori per i vari processi integrati nel BaBar-Software
- Basata su GEANT4.8
- I MC vengono prodotti inserendo le condizioni del detector per i vari periodi di presa dati (mensilmente)
 - Produzione di MC generici
 - Produzione di MC esclusivi
- Validazione basata su confronto con control-sample di dati
 - Attualmente l'accordo è a livello del 95%-99% in molti processi
- Successivamente gli eventi MC seguono la stessa validazione usata nei dati

EvtGen	B physics
KK2F	ISR + 2 fermion
Tauola	tau physics
JETSET	AKA Pythia, inclusive production of B and udsc
AfkQed	ISR + leptons or hadrons

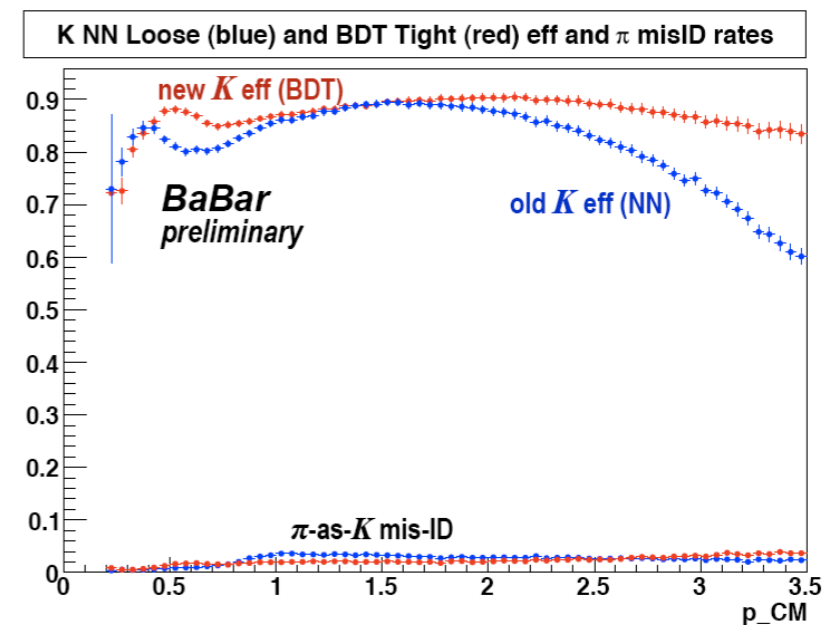
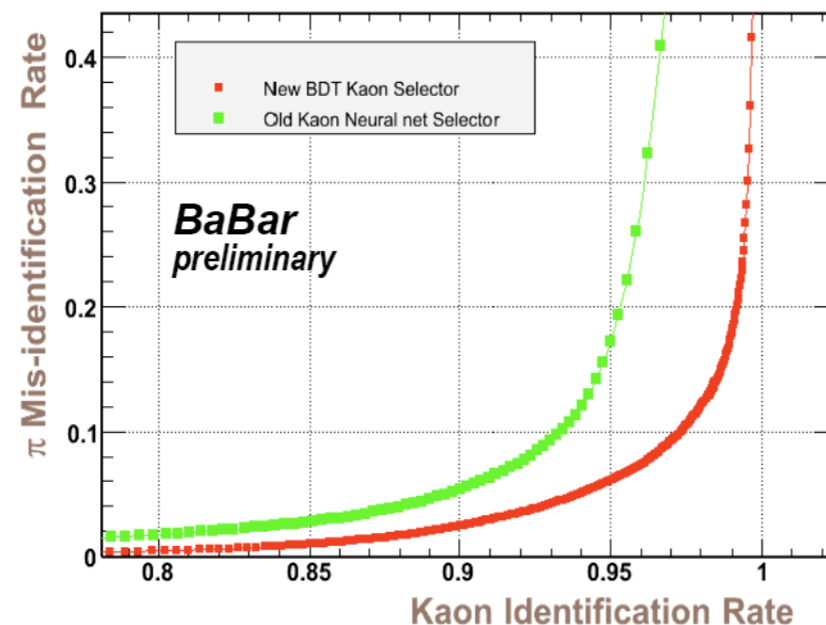


La release di ricostruzione

- La release contiene tutto il codice dell'esperimento
 - Vari pacchetti in CVS (~1000)
 - Ogni pacchetto ha almeno un coordinatore
 - Si distinguono i pacchetti per l'online, offline, analisi fisica
 - Un super-coordinatore per ognuno dei tre gruppi
- Le release vengono aggiornate continuamente, introducendo nuovi tag dei pacchetti
 - fast build delle release giornalmente
 - full build circa una volta a settimana
- Test di validazione per scegliere la release (non necessariamente la stessa) per
 1. acquisizione
 2. ricostruzione dei dati
 3. simulazione
 4. riprocessamento dati
 5. skimming
 6. analisi fisica

Riproccessamento dei dati

- Necessario quando la release incorpora significanti miglioramenti
 - Ultimo riproccessamento effettuato nel 2007, migliorando le efficienza di ricostruzione delle tracce
 - Aumento di efficienze nei segnali $\sim 10\%$, diminuzione dei fondi $\sim 10\%$
 - Impiegati circa 8 mesi in vari siti (SLAC, RAL, Padova)
- Attualmente è in corso il riproccessamento finale (The Final Game), atteso il completamento per Dicembre 2008
 - Migliorato tutto il codice per il PID, usando Boosted Decision Tree
 - Attualmente è in corso la validazione

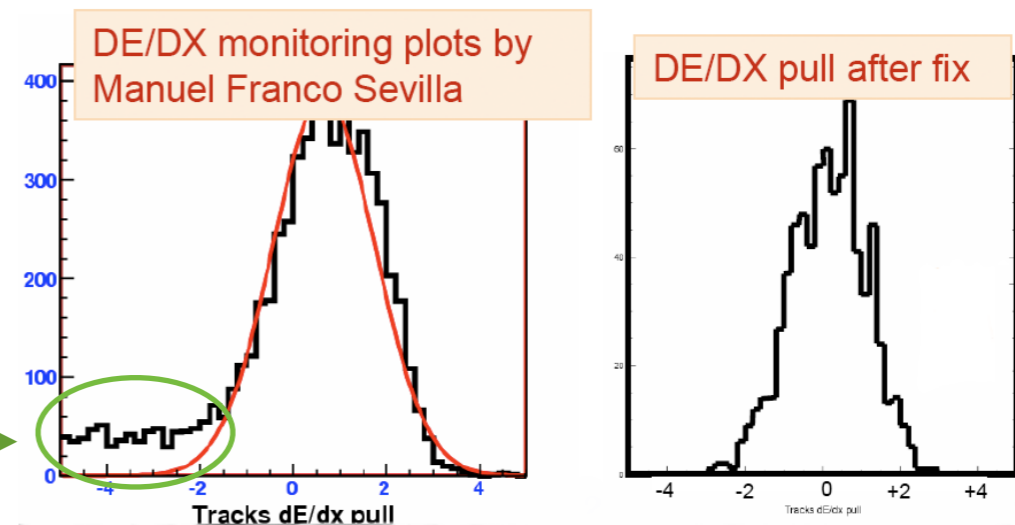


Data-quality

- I dati riprocessati vengono accuratamente validati, confrontando con dati ricostruiti con precedenti release
 - Analizzate varie variabili legate ad ogni sub-detector

- Un esperto per ogni sub-detector settimanalmente fornisce feedback sui vari dati riprocessati

Problema trovato nel PID (2007) →



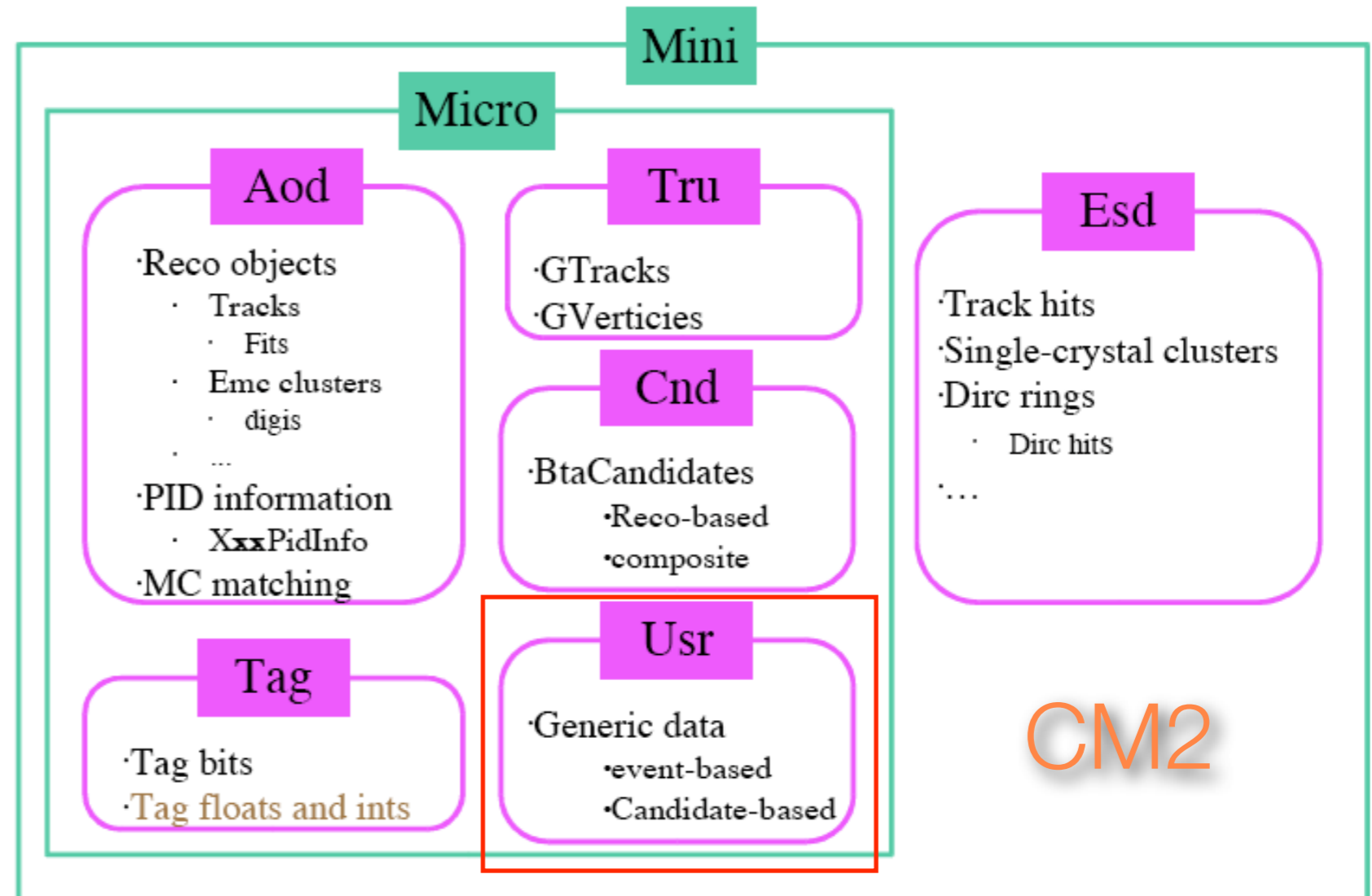
- Analizzate alcune variabili di fisica: numero di fotoni, carichi, quantità legate alle risonanze,..., ma non analisi complete
 - Analisi campione complete di particolari processi ($B \rightarrow D^* D^*$, $B \rightarrow K^* \gamma$), vengono eseguite su un sample ridotto di dati (**validazione 10 fb^{-1}**)
 - Solo per nuovi riprocessamenti si effettua una validazione su un più ampio dataset (**validazione 60 fb^{-1}**) per un gruppo maggiore di analisi

Skims

- I dati che superano la validazione vengono **raccolti in tag** per l'analisi fisica successiva
- Normalmente ogni analisi necessita di un set limitato di dati, ~3% del campione totale
 - vengono costruiti dei filtri per una preselezione degli eventi
 - Esempio analisi $B^0 \rightarrow \eta' K$: si chiede che nell'evento siano presenti almeno una η' e una K per formare un B con un certo taglio sull'impulso nel CM
 - ad ogni evento vengono assegnate le uscite dei **tagbit** che corrispondono all'uscita dei filtri
 - I dati con un certo set di tagbit vengono raccolti in stream (copie fisiche o copie virtuali dei dati)
 - Le stream vengono rese disponibili nei vari siti per l'analisi fisica: SLAC, RAL, IN2PR, CNAF, GRIDKA, ...

Organizzazione delle informazioni nei dati

- Standard file ROOT, con tutte le informazioni (Mini), tranne i Trigger Filter
- Gli analisti possono inserire le loro informazioni di fisica direttamente durante lo skimming dei dati (Usr Data)
 - Rapida ricostruzione per l'analista
- È possibile un riprocessamento dei dati direttamente da mini (mini-to-mini conversion) se non è richiesto di cambiare i Trigger Filter
- Diverse centinaia di skim prodotte



Physics Analysis

- La release per la ricostruzione ed analisi degli eventi usata dagli analisti
- Vengono verificati i pacchetti per
 - Tracking and Neutral reconstruction
 - PID
 - Particle Composition
 - Vertexing
 - *B*-Tagging
 - Fitting tools
- I dati vengono selezionati con tagli preliminari
- Le variabili di interesse per l'analisi (impulsi, posizioni, masse, energie, PID, ...) vengono salvate in root-tree
- La ricostruzione si spezza in diversi jobs che runnano parallelamente nella farm, ognuno producendo un root-tree di uscita

Root-tree esempio

 BRecMode	 KsPionSvtPattern22	 energyLER	 gamSecondMoment3
 BVertexChi2	 KsPionnDof12	 etaCharge1	 gamStatus3
 BVertexCovZ	 KsPionnDof22	 etaCosHel1	 gamTheta3
 BVertexNDof	 KsPionprobChi212	 etaEnergy1	 lowerID
 BVertexProb	 KsPionprobChi222	 etaEnergyCMS1	 mES
 BVertexStatus	 KsPionstartFoundRange12	 etaMass1	 mES_Err
 BVertexZ	 KsPionstartFoundRange22	 etaMassReco1	 mHat
 BcategoryTag04	 KsPiontrackLength12	 etaMomentum1	 mHat_Err
 Bcharge	 KsPiontrackLength22	 etaMomentumCMS1	 mMiss
 BtagTag04	 KsPx2	 etaPi0Veto1	 mMiss_Err
 DE	 KsPy2	 etaPx1	 mcBRecID
 DE_Err	 KsPz2	 etaPy1	 mcBRecMode
 KsCtau2	 KsVertexProb2	 etaPz1	 mcBRecZ
 KsCtauErr2	 Leg0	 evtNum	 mcBTagID
 KsDeltaVtxIP2	 Leg2	 gamDistGams3	 mcBTagMode
 KsDistanceXY2	 R2all	 gamDistTrks3	 mcBTagZ
 KsEnergy2	 bEnergy	 gamEnergy11	 mcIsCharmed
 KsFlightLength2	 bMass	 gamEnergy21	 mcMatchInfoDiff
 KsFlightLengthXY2	 bMass_Err	 gamEnergy3	 mcMatchInfoNoMatch
 KsKinkAngle2	 bPx	 gamEnergyB3	 mcMatchMode
 KsMass2	 bPy	 gamEnergyCMS3	 mcNBDauGamma0
 KsMomentum2	 bPz	 gamLat3	 mcNBDauGamma1
 KsMomentumCMS2	 combMassReco10	 gamNBumps3	 mcNBDauK0
 KsPionChi212	 combMassReco20	 gamNXtal3	 mcNBDauK1
 KsPionChi222	 combMassReco21	 gamPhi3	 mcNBDauKL0
 KsPionNDchHits12	 cosBBeam	 gamPi0EtaVeto3	 mcNBDauKL1
 KsPionNDchHits22	 cosBTBeam	 gamPxCMS3	 mcNBDauMu0
 KsPionNSvtHits12	 cosBTROE	 gamPyCMS3	 mcNBDauMu1
 KsPionNSvtHits22	 energyCMS	 gamPzCMS3	 mcNBDauOther0
 KsPionSvtPattern12	 energyHER	 gamR2prime3	 mcNBDauOther1

I dati a Milano: Selezione degli eventi

- I root-tree prodotti vengono copiati a Milano
- Concateniamo i vari file in **catene** (classe **TChain**)
- Per la selezione finale degli eventi usiamo un **selettore** (classe **TSelector**)
 - Abbiamo **esteso le funzionalità** della classe TSelector
 - Riutilizzare il codice più possibile
 - Tutto (o quasi) è configurabile senza modificare il codice, tramite opzioni
 - Ogni analisi ha il suo proprio selettore con specifici tagli
- L'ottimizzazione dei tagli, il calcolo delle efficienze, la stima del fondo sono tutte implementate all'interno del selettore
- L'uscita del selettore è interfacciata con un foglio di calcolo (XML) per calcoli successivi e consente di scrivere direttamente tabelle LaTeX (Pietro B.)
- Gli eventi che superano i tagli del selettore vengono salvate in un root-tree per l'input al fit Maximum Likelihood per l'estrazione dei risultati finali

Stima del fondo

- Ricostruzione degli eventi MC di fondo generici
- Salviamo le informazioni relative alla MC generazione di tutti gli eventi ricostruiti
 - Usiamo questa informazione per capire che canali danno fondo dopo la selezione con il selettore
- Studi più approfonditi sono successivamente fatti con MC esclusi

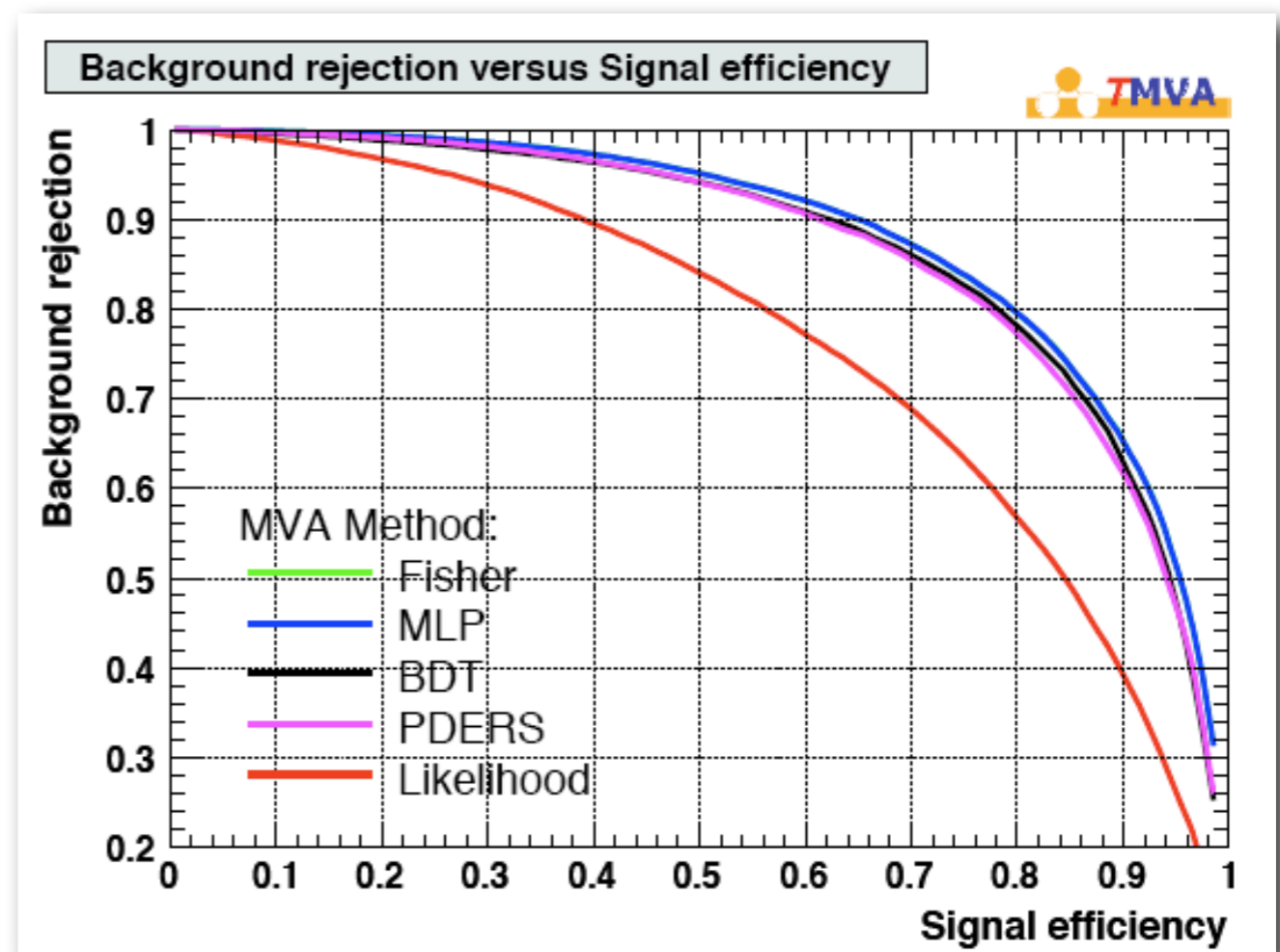
Name	Title
[anti-B0 -> rho- D+]	185
[B0 -> rho+ D-]	177
[anti-B0 -> rho- D*+]	149
[B0 -> rho+ D*-]	138
[B0 -> a_1+ D-]	102
[anti-B0 -> D*+ a_1-]	100
[anti-B0 -> D*+ pi-]	98
[B0 -> D*- pi+]	93
[anti-B0 -> a_1- D+]	83
[B0 -> D*- a_1+]	75
[B0 -> D*- nu_mu mu+]	65
[B0 -> D- pi+]	62
[anti-B0 -> D*+ anti-nu_mu mu-]	60
[anti-B0 -> D+ pi-]	47
[anti-B0 -> a_10 omega]	44
[B0 -> a_10 omega]	42
[anti-B0 -> omega rho- pi+]	38
[anti-B0 -> omega D0]	36
[B0 -> omega rho- pi+]	36
[anti-B0 -> rho0 D0]	34

Soppressione del fondo

- Analisi Multivariate impiegando varie tecniche
 - Ottimizzazione dei tagli Multivariata
 - Discriminante di Fisher
 - Reti Neurali
 - Bagging & Boosted Decision Tree
- Pacchetti usati:
 - **TMVA -- Toolkit for Multivariate Data Analysis:** General framework for Multivariate data analysis
<http://tmva.sourceforge.net/>
 - **SPR -- StatPatternRecognition:** General framework for Multivariate data analysis
<http://sourceforge.net/projects/statpatrec>
- Tecniche molto potenti per la separazione segnale/fondo, comunemente usate nel gruppo di Milano (e in BaBar in generale)

Tecniche avanzate di soppressione segnale/fondo

- Recentemente ho dato un talk alla iCSC08 al CERN sull'argomento <http://csc.web.cern.ch/CSC/2008/iCSC2008/iCSC2008.htm>
- Riferimento bibliografico:
Trevor Hastie, Robert Tibshirani, Jerome Friedman, "The Elements of Statistical Learning - Data Mining, Interference, and Prediction", Springer Series in Statistics (2003)
- I software sono facili da utilizzare, ben documentati, integrati in ROOT
 - Facile fare confronti tra le varie tecniche

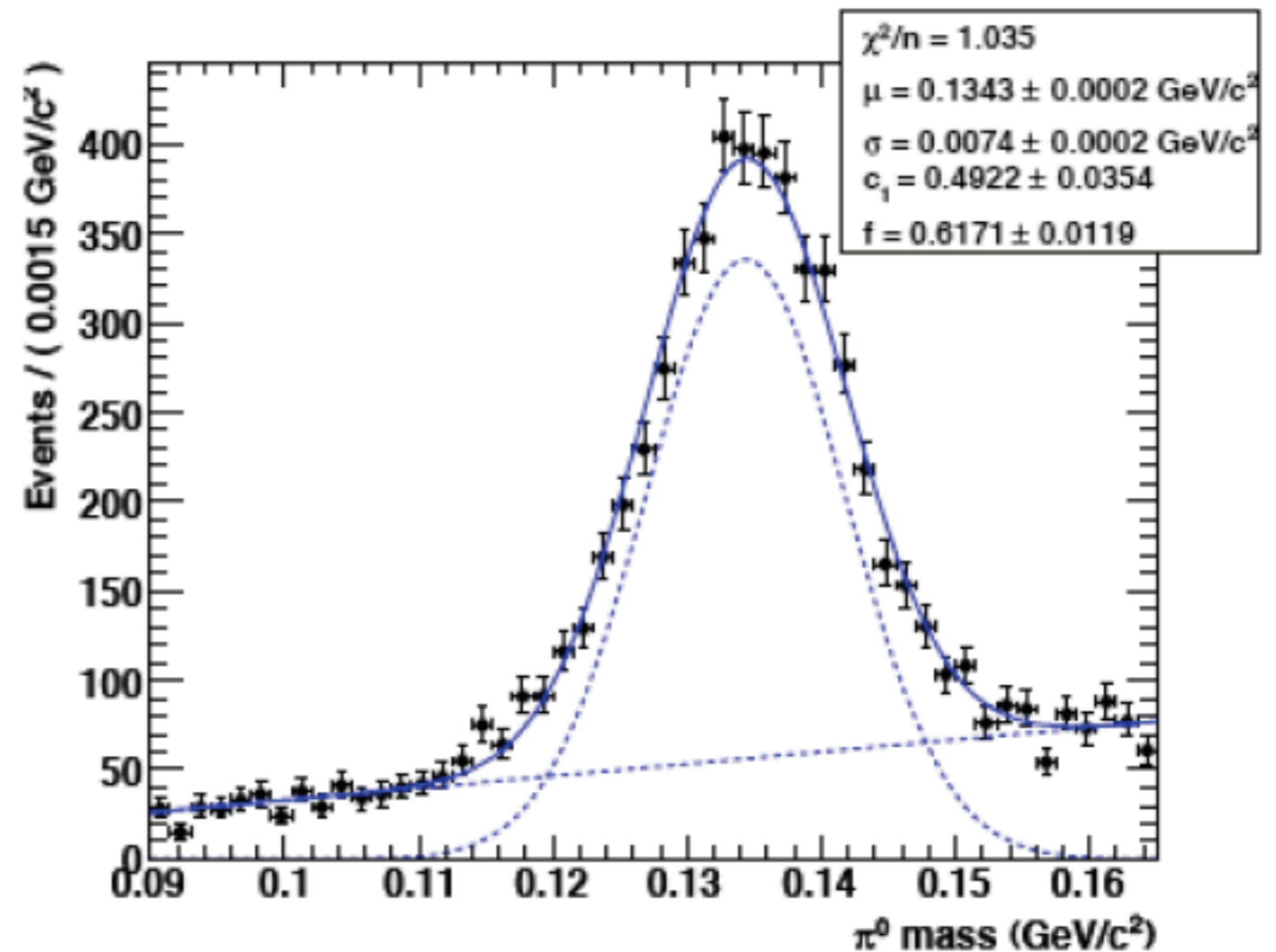


Fit Maximum Likelihood

- Fit Multivariati Unbinned basato sul pacchetto RooFit
 - Integrato in ROOT
 - Esperienza con RooFit fin dagli inizi del suo sviluppo (entrambi gli autori, D. Kirkby e W. Verkerke, sono/erano membri di BaBar)
 - Praticamente usato in tutte le analisi di BaBar: grande esperienza accumulata nell'esperimento
 - Recente workshop tenuto a SLAC a Dicembre
http://www.slac.stanford.edu/BFROOT/www/doc/Workshops/2007/BaBar_RooFit/Agenda.html
- Comunque RooFit è pur sempre un insieme di classi C++ e spetta all'utente dover "lavorare" con il codice per scriversi il proprio programma
 - **Poco versatile**: se le analisi hanno aspetti in comune, è possibile pensare ad un sistema di scripting per il "riciclo" del codice
 - **Basso livello**: RooFit fornisce un set di classi base, non sufficienti per effettuare proiezioni, NLL scan, Upper Limit, esperimenti MC simulati, ...
- Serve qualcosa di più generale che permetta di configurare ROOT & RooFit

MiFit

- Interfaccia ROOT & RooFit
- Interamente configurabile da semplici file di testo
 - Non serve alcuna conoscenza approfondita del codice
 - Costrutti superiori per la dichiarazione delle varie operazioni
- Generalizzato a qualsiasi tipo di analisi ML unbinned
- Attualmente il collo di bottiglia di molte analisi in BaBar è diventato il tempo di esecuzione dei fit: analisi Dalitz-plot a molti parametri liberi e molti eventi impiegano ore (giorni)
 - Diventa molto interessante usare tecniche di calcolo parallelo, attualmente mio oggetto di studio



Coclusione

- L'esperienza di BaBar dice che è fondamentale avere un controllo attento dei dati e di tutto il codice di ricostruzione e di analisi fisica
- Particolarmente importante quando si ha un esperimento concorrente: Belle
- Tecniche avanzate consentono a BaBar di restare competitivo (stessa significatività) con Belle malgrado loro abbiano l'80% in più dei dati rispetto all'attuale campione di BaBar
 - particolarmente importante considerato che BaBar smetta la presa dati alla fine di questo mese, mentre Belle prosegue per almeno altri 2 anni

Fine del talk



A
T
L
A
S