

Pattern Recognition and Event Reconstruction in Particle Physics Experiments

R. Mankel¹

Deutsches Elektronen-Synchrotron DESY, Hamburg

Abstract

This report reviews methods of pattern recognition and event reconstruction used in modern high energy physics experiments. After a brief introduction into general concepts of particle detectors and statistical evaluation, different approaches in global and local methods of track pattern recognition are reviewed with their typical strengths and shortcomings. The emphasis is then moved to methods which estimate the particle properties from the signals which pattern recognition has associated. Finally, the global reconstruction of the event is briefly addressed.

¹Email: Rainer.Mankel@desy.de

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 3 |
| 2 | Basics | 5 |
| 2.1 | Detector Layouts | 5 |
| 2.1.1 | Forward or fixed target geometry | 6 |
| 2.1.2 | Collider detector geometry | 7 |
| 2.2 | Typical Tracking Devices | 9 |
| 2.2.1 | Linear single-coordinate measurements | 9 |
| 2.2.2 | Radial single-coordinate measurements | 10 |
| 2.2.3 | Stereo angles | 12 |
| 2.2.4 | Three-dimensional measurements | 13 |
| 2.3 | Track Models and Parameter Representations | 14 |
| 2.3.1 | Forward geometry | 14 |
| 2.3.2 | Cylindrical geometry | 15 |
| 2.4 | Parameter Estimation | 15 |
| 2.4.1 | Least squares estimation | 16 |
| 2.4.2 | The Kalman filter technique | 17 |
| 2.5 | Evaluation of Performance | 19 |
| 2.5.1 | The reference set | 20 |
| 2.5.2 | Track finding efficiency | 20 |
| 2.5.3 | Ghosts | 21 |
| 2.5.4 | Clones | 21 |
| 2.5.5 | Parameter resolution | 22 |
| 2.5.6 | Interplay | 22 |
| 3 | Global Methods of Pattern Recognition | 23 |
| 3.1 | Template Matching | 23 |
| 3.2 | The Fuzzy Radon Transform | 25 |
| 3.3 | Histogramming | 26 |
| 3.4 | Neural Network Techniques | 33 |
| 3.4.1 | The Hopfield neuron | 34 |
| 3.4.2 | The Denby-Peterson method | 35 |
| 3.4.3 | Elastic arms and deformable templates | 40 |
| 4 | Local Methods of Pattern Recognition | 50 |
| 4.1 | Seeds | 50 |
| 4.2 | 2D Versus 3D propagation | 51 |
| 4.3 | Naïve Track Following | 53 |
| 4.4 | Combinatorial Track Following | 54 |
| 4.5 | Use of The Kalman Filter | 55 |
| 4.6 | Arbitration | 55 |

| | | |
|----------|--|-----------|
| 4.7 | An Example for Arbitrated Track Following | 56 |
| 4.7.1 | Algorithm | 56 |
| 4.7.2 | Parameters | 57 |
| 4.8 | Track Following And Impact of Detector Design Parameters . . . | 58 |
| 4.8.1 | Influence of detector efficiency | 59 |
| 4.8.2 | Effect of detector resolution | 59 |
| 4.8.3 | Influence of double track separation | 59 |
| 4.8.4 | Execution speed | 63 |
| 4.9 | Track Propagation in a Magnetic Field | 64 |
| 5 | Fitting of Particle Trajectories | 66 |
| 5.1 | Random Perturbations | 66 |
| 5.2 | Treatment of Multiple Scattering | 66 |
| 5.2.1 | Impact parameter and angular resolutions | 71 |
| 5.2.2 | Momentum resolution | 71 |
| 5.2.3 | Effects of fit non-linearity | 73 |
| 5.2.4 | Contributions of different parts of the spectrometer | 74 |
| 5.2.5 | Parameter covariance matrix estimation | 75 |
| 5.2.6 | Goodness of fit | 75 |
| 5.3 | Treatment of Ionization Energy Loss And Radiation | 79 |
| 5.3.1 | Ionisation energy loss | 79 |
| 5.3.2 | Radiative energy loss | 79 |
| 5.3.3 | Radiation energy loss correction within the magnetic field | 85 |
| 5.4 | Robust Estimation | 88 |
| 6 | Event Reconstruction | 89 |
| 6.1 | Vertex Pattern Recognition | 89 |
| 6.2 | Vertex Fitting | 91 |
| 6.3 | Kinematical Constraints | 92 |
| 7 | Concluding Remarks | 93 |
| | Acknowledgement | 93 |
| | References | 93 |

1 Introduction

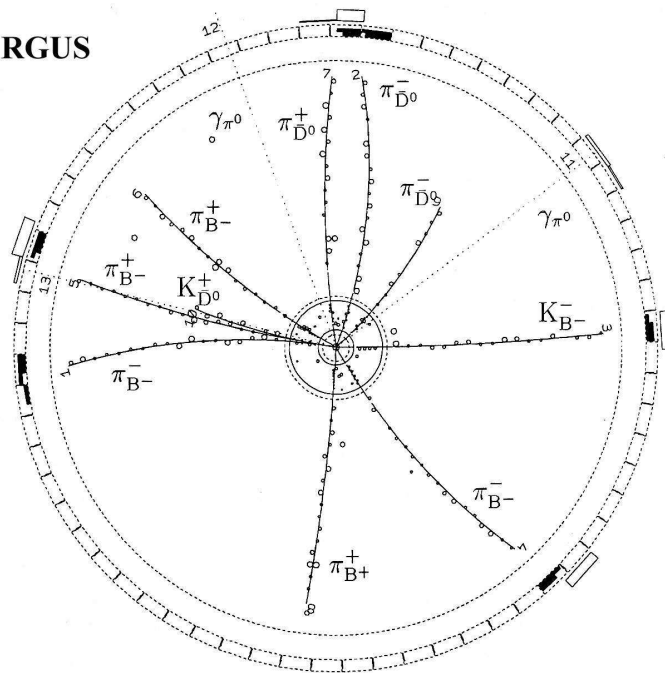
Scientific discovery in elementary particle physics is largely driven by the quest for higher and higher energies, which allow delving ever more deeply into the fine structure of the microscopic universe. Higher energies lead in general to an increased multiplicity of particles. Since the acceleration of electrons is limited either by synchrotron radiation in case of storage rings, or by field gradients in case of linear colliders, multi-TeV energies are in the near future only accessible by accelerating hadrons, the collision of which generates even more particles.

Reconstruction of charged particles from signals of tracking detectors in spectrometers has always shown aspects of a discipline of art, since the variety of experimental setups lead to development of very diverse pattern recognition methods, which could not easily be ranked among each other. An general overview has been given in an earlier review [1]. It is remarkable that even today, no generally accepted standard software package exists which performs track finding in a variety of detector setups, a situation which is in marked contrast e.g. to detector simulation. A new generation of experiments is now emerging in which the track density is so high that success will crucially depend on the power of the reconstruction methods. One example for the development in tracking demands over 15 years is illustrated in fig. 1, which shows in direct comparison an event from the experiment ARGUS [2], which took data of e^+e^- collisions in the Υ range in the period 1982–1992, and the ATLAS experiment [3] which is currently under construction and will operate from 2007 on with proton collisions at the LHC. The new experiments also require huge computing resources for reconstruction of their data. Since track finding is usually the most time consuming part in reconstruction, the sophistication and economy of pattern recognition methods has considerable impact on the computing effort.

Pattern recognition plays an important rôle also in other detector components, for example cluster reconstruction in calorimeters, or ring finding in ring imaging Čerenkov detectors (RICH). It is however in track reconstruction where the new generations of experiments pose the most crucial challenges. This article will therefore focus on track reconstruction as well as to related aspects of event reconstruction.

The first of the following chapters will provide an introduction into basic detector concepts and tracking devices and summarize mathematical tools for estimating parameters and performance that will be used later on. The two following chapters focus on track pattern recognition with various methods, including applications in several experiments. The next chapter then concentrates on parameter estimation from particle trajectories, which is – in contrast to track finding – in principle a straight-forward mathematical problem, but contains several detailed issues worth mentioning. The last chapter briefly discusses some track-related aspects of event reconstruction.

ARGUS



ATLAS

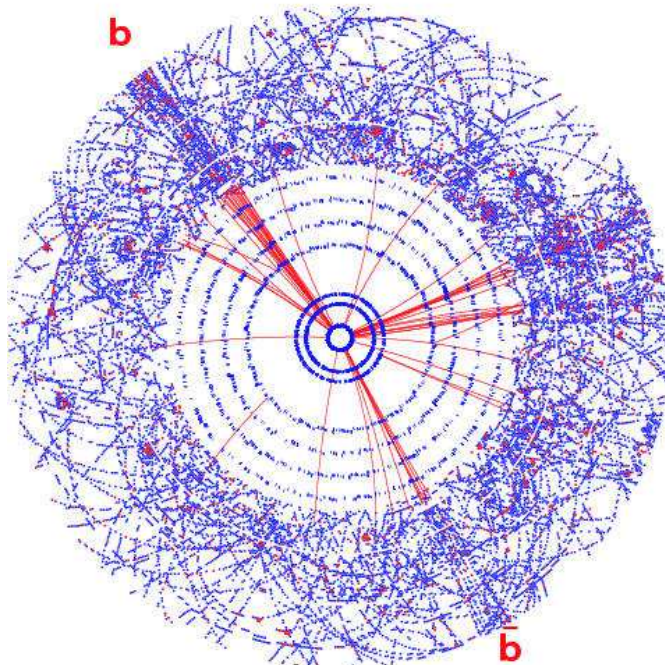


Figure 1: Comparison of event complexity in the experiments ARGUS and ATLAS. The ARGUS event (top) consists of two reconstructed B mesons, one of them being a candidate for the charmless decay $B^- \rightarrow K^- 4\pi^\pm$ (from [2]). The ATLAS display (bottom) shows a simulation of an event in the inner detector with a Higgs boson in the decay mode $H^0 \rightarrow b\bar{b}$, including the pileup at full LHC luminosity (from [3]).

2 Basics

This section provides a brief introduction into the basic elements influencing event reconstruction. It is not intended to cover the subject of particle detectors in full detail, instead the detector literature (see for example [4, 5, 6]) is referred to.

2.1 Detector Layouts

Modern detectors in high energy physics are usually sampling detectors. The detector volume is filled with devices which the particles traverse and in which they leave elementary pieces of information, as e.g. an excitation in a solid-state detector, a primary ionization in a gaseous chamber or an energy deposition in a sensitive volume of a calorimeter. The event record of an experiment consists of the amassed volume of the signals from all particles of an interaction – or possibly even several interactions – joined together. After sorting out which bits of information are related to the same particle – this process is called *pattern recognition* – the kinematical properties of each particle have to be reconstructed, to reveal the physical nature of the whole event.

In general, experiments nowadays strive to record the interaction as a whole, with all (significant) particles produced in the process. This has lead to the development of 4π detectors, where almost the whole solid angle region, as seen from the interaction, is covered.

In general, two main concepts have to be distinguished, which will be discussed in the following.

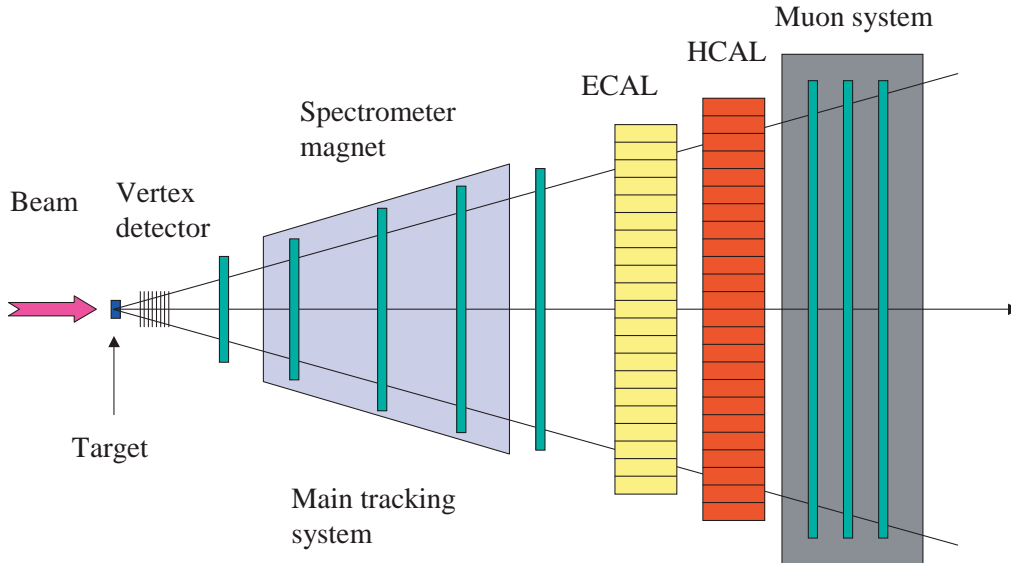


Figure 2: Typical geometry of a forward spectrometer, as used e.g. in fixed-target setups.

2.1.1 Forward or fixed target geometry

When the interaction is generated by an incident beam hitting a fixed target, the centre-of-mass system of the participating particles is seen under a strong Lorentz boost, and the emerging particles are moving within a cone into the forward direction. In this case, the detector setup must cover this forward cone with instrumentation, while the more backward part of the solid angle is generally neglected. This scenario is called a *forward detector geometry*. Similar situations exist where the dynamics of the interaction result in all relevant particles to be produced under a huge Lorentz boost, like heavy flavour production at large hadron colliders.

Figure 2 schematically shows a forward detector geometry as it is used in fixed target experiments. The event is generated through collision of a beam particle with a nucleus in the target. Because of the momentum of the incident beam particles, the whole event is seen under a Lorentz boost in the beam direction, so that the emerging particles are confined to a cone whose opening angle depends on the typical transverse momenta generated in the interaction, and the size of the Lorentz boost.

The main components of a typical forward spectrometer are:

- the vertex detector, which is a precision tracking system very close to the interaction point. Its main purpose is the improvement of track resolution near the interaction point which allows reconstruction of secondary vertices or distinction of detached tracks which is used e.g. for the tagging of heavy flavour decays.
- the spectrometer magnet with the main tracking system, which measures trajectories of charged particles and determines their momentum and charge sign from the curvature.
- the calorimeter system, which is often split into an electromagnetic and a hadronic part. The calorimeter allows identification of electrons and hadrons by their deposited shower energy, and very often provide essential signals for the trigger system. The calorimeter can also measure energies of individual neutral particles, in particular photons, though the actual capability in this task depends strongly on the particle density in the event.
- the muon detector, which consists of tracking devices in combination with absorbers. Only muons are able to traverse the intermediate material, and are then measured in the dedicated tracking layers.

The design of a forward spectrometer is influenced by several factors. The sheer size of the tracking volume depends on the leverage required for the momentum resolution, since at sufficiently high momentum the resolution is inversely proportional to the integral of the magnetic field along the trajectory [7], as will

be discussed in more detail in sec. 5. Depending on the scope of the experiment, further detector components may be introduced to provide particle identification, for example ring-imaging Čerenkov counters (RICH) or transition radiation detectors (TRD).

2.1.2 Collider detector geometry

When two beams collide head-on, the centre-of-mass system of the interactions is either at rest or moving moderately. In this case, the detector should try to cover the full solid angle. This beam setup usually leads to cylindrical detector layouts with a solenoid field parallel to the beam axis (fig. 3). In comparison to

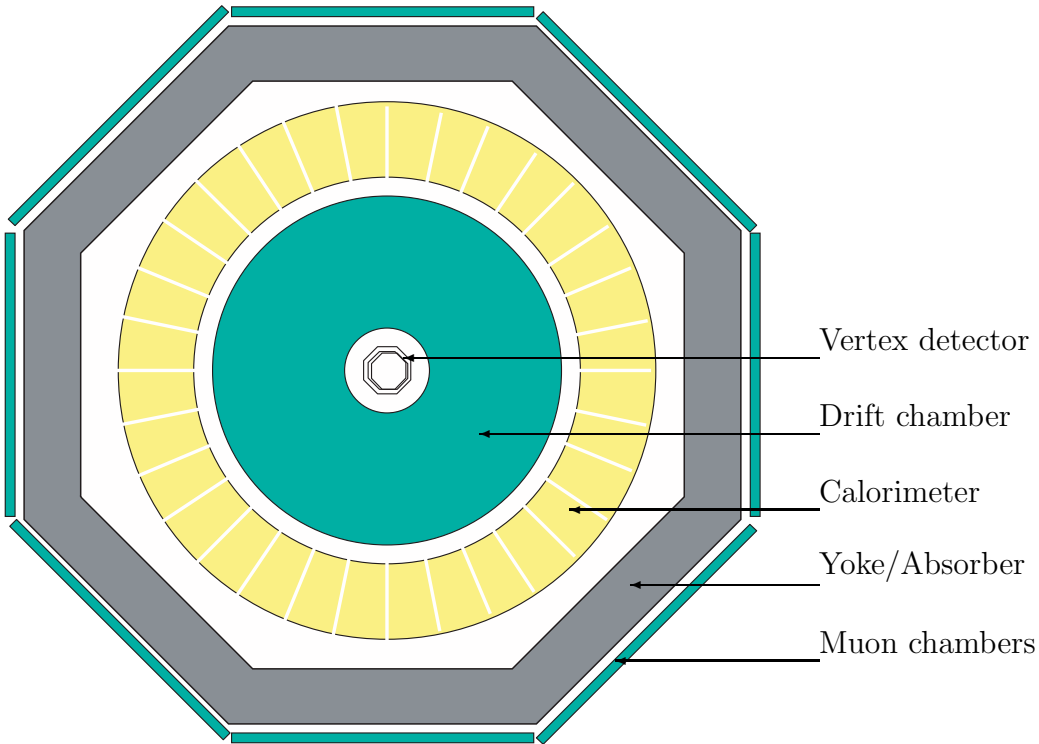


Figure 3: Typical setup of a collider detector.

the forward geometry detector, the cylindrical geometry differs in several details:

- the vertex detector requires modules parallel to the beam, at least in the central part of the angular acceptance, often referred to as the *barrel part*.
- the main tracking system is generally contained in the magnetic field. Coil and yoke of the magnet usually have to be within the detector volume, where the general choice is to have the coil between drift chamber and calorimeter, where particles traverse it before their energy being measured in the calorimeter, or to make it large enough to enclose the calorimeter,

which may be more costly to build and operate and where the field may have adverse effects on the calorimeter itself.

- the calorimeter system now requires barrel and end cap parts to cover the solid angle. A main functionality at high energy colliders is the measurement of jets.
- for the muon detector, the yoke of the solenoid lends itself readily as absorber.

2.2 Typical Tracking Devices

2.2.1 Linear single-coordinate measurements

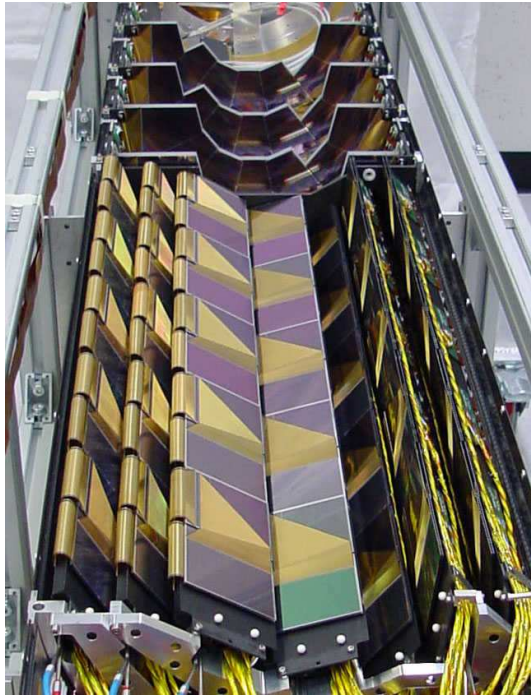


Figure 4: Lower half barrel part of the Zeus micro-vertex detector

A widespread type of tracking device measures one coordinate of the particle whose trajectory intersects the device. A good example for this type represent silicon strip detectors, which are semiconductor-based devices structured in strips typically down to widths of $25\text{ }\mu\text{m}$. Each strip works like a small diode, with a voltage applied such that the border area is depleted and the resistance is high. A traversing charged particle will then create pairs of electrons and corresponding holes which drift apart under the voltage and can be registered as a pulse. In general several strips will register a signal under traversal of a particle, and the pulse heights of the participating channels can be evaluated with suitable clustering algorithms, for example centre-of-gravity based, and determine the location at which the particle has passed. Solid-state detectors are presently the tracking devices with the highest spatial resolution, and they are often installed very close to the interaction region as *vertex detectors* where they allow or improve the reconstruction of primary and secondary vertices. Another favourable property of solid-state detectors is their resilience against radiation damage. The current limitation is in the size of individual detector modules, which makes them expensive for coverage of large volumes. Figure 4 shows the micro-vertex detector of

the ZEUS experiment [8], prior to its installation in 2001.

2.2.2 Radial single-coordinate measurements

The size of the tracking volumes is important, since momentum measurement requires the particle to traverse a magnetic field, where the length of the path provides the leverage that determines the precision of the momentum reconstruction. This is one of the reasons why gaseous chambers, in particular drift chambers are very commonly employed when large areas have to be covered.

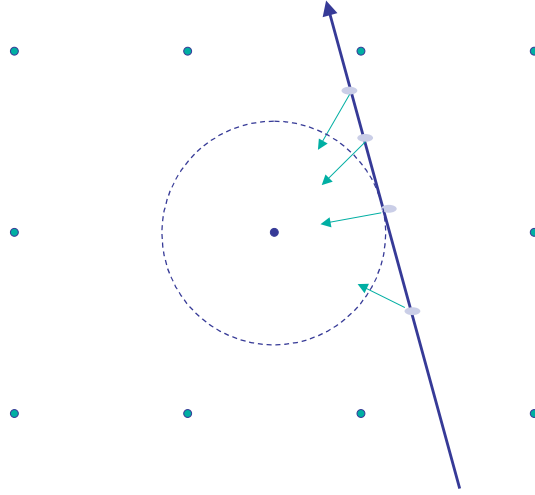


Figure 5: Schematic view of a drift chamber cell. The filled circles indicate wires, with the sense wire in the middle of the cell and the field wires on the outside. The black arrow shows the trajectory of a particle, the grey arrows denote primary ionization charges drifting towards the sense wire.

The basic principle of the drift chamber is displayed in fig. 5. A drift cell consists of an anode wire in the centre and an arrangement of field wires. The geometry shown is very similar to that in the ARGUS drift chamber [9] (see also fig. 34 in section 4). The drift cell need not be of rectangular shape, in the drift chamber of the BaBar experiment, for example, it is hexagonal [10]. Along the path of the particle, primary ionization occurs. The charges drift to the anode wire, where they create a locally confined avalanche of particles within the large electrical field close to the wire. This effect results in a multiplication of the ionization which is called *gas amplification*. The rising edge of the signal picked up by the anode wire triggers a time-to-digital converter (TDC) which then measures the time until a common stop signal. This allows measuring of the drift time for those charges that are the first to arrive. In the simplest case, the drift field will be shaped such that the drift velocity is uniform, and the time resolution can be directly transformed into a uniform resolution of the drift

distance. In practice, numerous effects can lead to a non-linear *drift-time/space* relation, and the spatial resolution will depend on the precise location of the traversal of the particle.

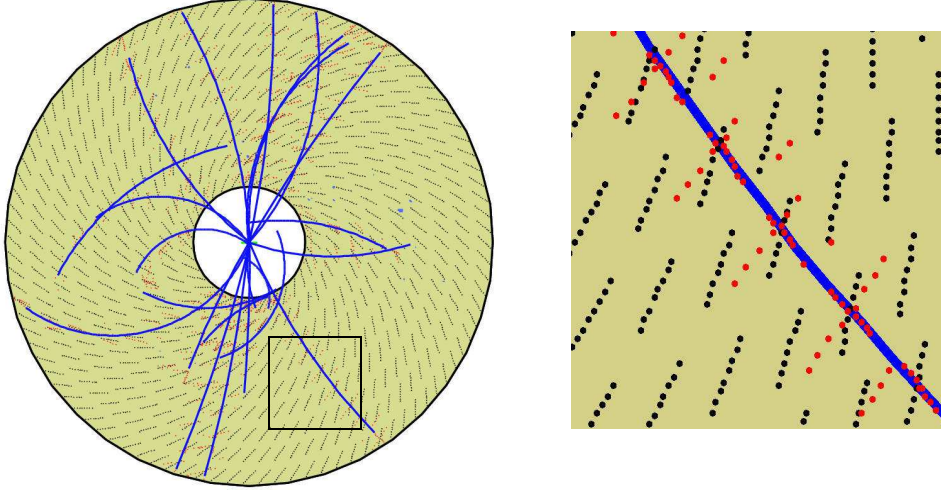


Figure 6: Left: event display from the ZEUS central tracking detector (CTD), showing sense wires and reconstructed tracks. Right: closeup around the track in the lower left area. The black dots represent the sense wires, the grey dots indicate the drift distance end points on both sides of the corresponding wire.

Since the time measured by the TDC corresponds to the arrival of the first charges, usually those with the smallest distance to the wire, the drift chamber measures the distance of closest approach of the particle to the wire. In cases where more than one particle traverses the same drift cell within the same interaction window, in general only the particle closest to the wire is registered. This effect may cause complications for pattern recognition which depend on the degree of occupancy. Another typical property of drift chambers is that the single measurement cannot distinguish on which side of the wire the particle has traversed; this uncertainty is called *left-right ambiguity*. In the worst case, left-right ambiguity may lead to a *mirror track* that cannot be distinguished from the real one. Concepts have therefore been developed how to design drift chambers such that left-right ambiguity can be resolved in all cases, e.g. the *butterfly geometry* [11].

Drift in gases is influenced also by magnetic fields. The deviation of the gas drift direction from the vector of the electric field is described by the *Lorentz angle*. Figure 6 shows an event display of the central tracking detector (CTD) of the ZEUS experiment, in the view along the beam axis, which has been created

using the tool described in [12]. The Lorentz angle in this case is 45° , and it is reflected in the design of the cell structure.

2.2.3 Stereo angles

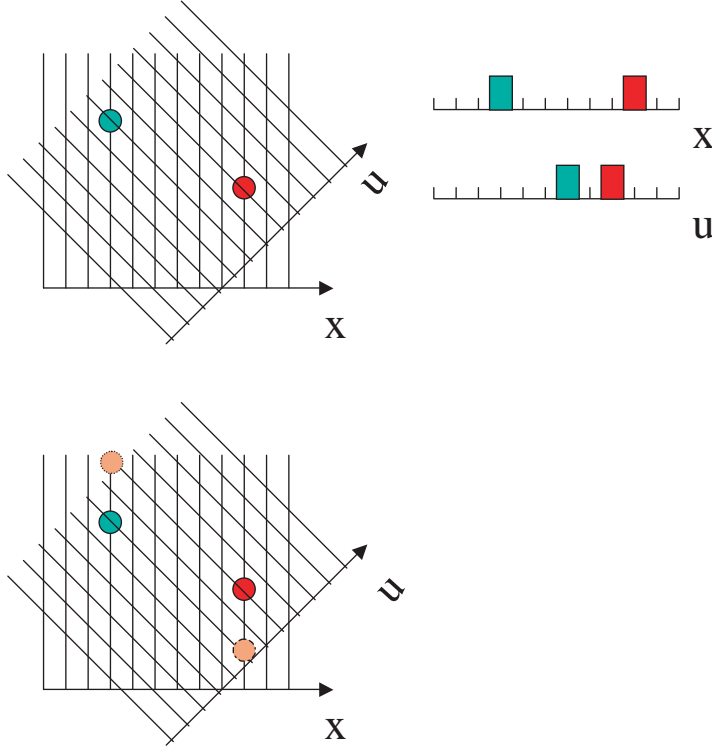


Figure 7: Hit ambiguities with two stereo views

Devices measuring single coordinates do not provide three-dimensional² points on a trajectory, but measure only in a projected space. While such devices can be very economic in the sense that a relatively small number of channels is needed to cover a region at good resolution, 3D information can only be obtained by combining several projections, usually named *stereo views*. While two views are in principle sufficient to reconstruct spatial information, the presence of more than one track leads in general to ambiguities regarding the assignment of projected information. This is illustrated in fig. 7, where two particles are measured in two strip detector views of 0° (x) and 45° (u). Ambiguity in the assignment of the measured hits in the x and u views to each other leads to the reconstruction of two ghost points. This illustrates that in general at least three views are necessary to avoid this kind of ambiguities. On the other hand, in special cases of limited

²The shorthands 2D (two-dimensional) and 3D (three-dimensional) will frequently be used in the following.

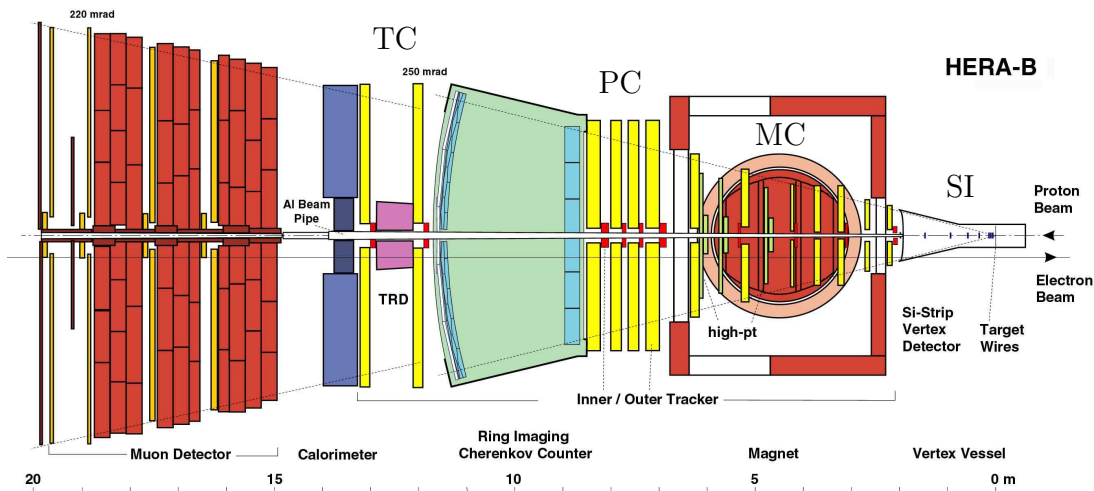


Figure 8: Layout of the HERA-B spectrometer. The labels TC, PC, MC and SI indicate groups of tracking stations that comprise the vertex and main tracking system.

track density, the use of only two views may be justified, since in this case the majority of ghosts may be discarded for geometrical reasons. This can already be guessed from fig. 7: since the true tracks are well separated, the uppermost ghost combination is already just outside the chamber acceptance of the u view. Such concepts are called *all-stereo* designs.

An example for a spectrometer that combines several types of single-coordinate measurements is the HERA-B detector [13, 14, 15] which is shown in fig. 8. The vertex detector (labelled *SI*) consists of eight superlayers of silicon strip detectors with four different stereo angles. The design of the main tracker is structured into the three areas within the magnet (*MC*), between magnet and RICH (*PC*) and between RICH and calorimeter (*TC*), it contains 13 superlayers of honeycomb drift chamber modules for the outer area and 10 superlayers of micro-strip gaseous chambers (MSGC) for the region close to the beam³.

2.2.4 Three-dimensional measurements

In general pattern recognition will benefit considerably if the tracking device itself is able to measure 3D space points. A modern example is solid-state pixel detectors, as for example the CCD-based vertex detector of the SLD experiment [17],

³The layout of tracking stations has been modified later with the shift of emphasis away from B physics.

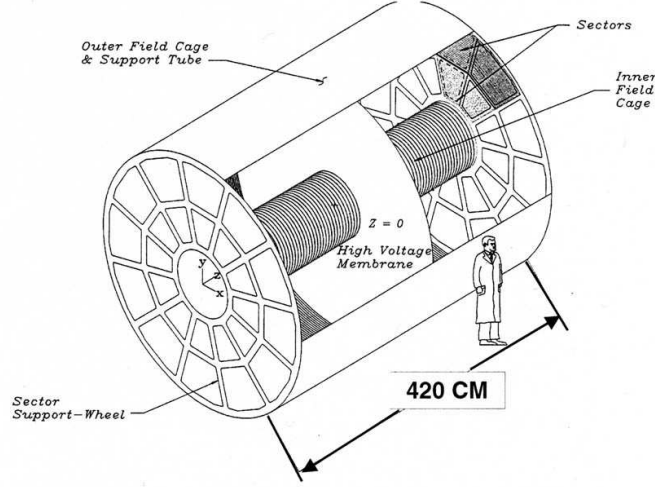


Figure 9: TPC of the STAR experiment (from [16]).

where the pixels have a size of $20 \times 20 \mu\text{m}^2$. A gaseous detector capable of covering large tracking volumes with 3D measurement is the *time projection chamber* (TPC). Figure 9 shows the TPC of the STAR experiment [16]. The gas volume itself is free of wires; instead, an axial electrical field, produced with the help of a membrane electrode in the middle plane, lets the primary charges drift to the anodes at the end caps, where they are registered, for example with multi-wire proportional chambers with pad readout. While this provides a direct measurement of the x and y coordinates, the z coordinate is inferred from the time measurement. The magnetic field is also axial, and plays an important rôle in limiting diffusion effects during the drift.

2.3 Track Models and Parameter Representations

2.3.1 Forward geometry

In the forward geometry, the interaction region lies very often in an area without magnetic field, since the spectrometer magnet is located further downstream. The natural choice of parameters, assuming that the z coordinate points down the spectrometer axis and x and y are the transverse coordinates, is then

x_0 the x coordinate at the reference z_0

y_0 the y coordinate at the reference z_0

$t_x = \tan \theta_x$ the track slope in the xz plane

$t_y = \tan \theta_y$ the track slope in the yz plane

Q/p the inverse particle momentum, signed according to charge

where z_0 denotes the location of a suitable reference plane transverse to the beam, for example at the position of the target, or at the nominal interaction point. The slope parameters allow for a convenient transformation of the parameters to a different reference z value, as is needed during vertex reconstruction. In cases of a very homogeneous magnetic field, it may be advantageous to substitute the parameter Q/p by Q/p_\perp , where p_\perp is the momentum in the plane transverse to the magnetic field, or by $\kappa = Q/R$, the signed inverse radius of curvature.

2.3.2 Cylindrical geometry

In collider detectors with cylindrical geometry, the magnetic field normally encompasses the whole tracking volume, including the interaction region where the particles are produced. In a homogeneous solenoid field, the particle trajectory will be a helix curling around an axis parallel to the magnetic field. Assuming the z coordinate is oriented along the detector axis, and the radius is given by $r = \sqrt{x^2 + y^2}$, typical track parameters given at a reference value $r = r_0$ may be

ϕ_0 the azimuth angle where the trajectory intersects the reference radius

z_0 the z value where the trajectory intersects the reference radius

ψ_0 the phase angle of the helix at the reference radius intersection, which corresponds to the angle of the tangent at this point

Q/R the signed inverse curvature radius of the helix

$\tan \lambda$ where $\lambda = \arctan p_z/p_\perp$ is the dip angle of the helix

2.4 Parameter Estimation

The estimation of the kinematical parameters of a particle, as position (or impact parameter), direction of flight and momentum at its point of origin from spatial measurements along its trajectory is generally referred to as track fitting. We assume at this point that the measurements related to a particle have been correctly identified in the pattern recognition step (which will be discussed in more detail in sections 3 and 4). A very general approach to parameter estimation is the *maximum likelihood method*, which shall not be discussed here in detail; instead we refer to the textbook literature [18, 19, 20, 21, 22]. The maximum likelihood method can take very general distributions of the observed variables into account, for example exponential distributions as they may occur when decay lengths are measured. However, its application in multi-parameter problems can be very complex, in particular the error analysis. In cases where the distribution of the random variables is Gaussian, at least approximately, the *least squares method* is

generally successful. Since many observables in track reconstruction do at least approximately follow a Gaussian distribution, this method will be focussed on in the following.

2.4.1 Least squares estimation

If the trajectory of a particle can be described by a closed expression $f_{\vec{\lambda}}(\ell)$, where $\vec{\lambda}$ stands for the set of parameters, ℓ is the flight path and f is the coordinate which could be measured, a set of measurements $\{m_i\}$ with errors $\{\sigma_i\}$ will provide an estimate of the parameters according to the least squares principle

$$X^2 = \sum \frac{(m_i - f_{\vec{\lambda}}(\ell_i))^2}{\sigma_i^2} \stackrel{!}{=} \min \quad (1)$$

One can easily convince oneself that in the case of normally distributed measurements m_i , the above impression is proportional to the negative logarithm of the corresponding likelihood function, which shows directly the equivalence of least squares principle and maximum likelihood principle for this case.

Symbolizing the derivative matrix⁴ of f with respect to the parameters as \mathbf{F} and the (diagonal) error matrix of the measurements as $\mathbf{V} = \text{diag}\{\sigma_i^2\}$, the expression to be minimized is

$$(\vec{m} - F\vec{\lambda})^T V^{-1} (\vec{m} - F\vec{\lambda}) \quad (2)$$

and requiring the derivative to vanish at the minimum leads to the matrix equation

$$F^T V^{-1} \vec{f} = F^T V^{-1} \vec{m} \quad (3)$$

In case of a linear problem, $\vec{f} = F\vec{\lambda}$, the above condition can be directly inverted

$$\vec{\lambda} = (F^T V^{-1} F)^{-1} F^T V^{-1} \vec{m} \quad (4)$$

and the estimated parameters are a linear function of the measurements. The matrix $(F^T V^{-1} F)^{-1}$ that needs to be inverted is of the shape $N_{\lambda} \times N_{\lambda}$ (where N_{λ} is the number of parameters describing the particle), which is inexpensive in terms of computation. Also the covariance matrix of the parameter estimate can be directly determined as

$$\text{cov}(\vec{\lambda}) = C_{\lambda} = (F^T V^{-1} F)^{-1} \quad (5)$$

The popularity of the least squares method can be attributed to its optimality properties in the linear case:

⁴We denote the derivative matrix as $\frac{\partial f}{\partial \lambda}$, where $\left(\frac{\partial f}{\partial \lambda}\right)_{ij} = \frac{\partial f_{\vec{\lambda}}(\ell_i)}{\partial \lambda_j}$.

- the estimate is unbiased, i.e. the expectation value of the estimate is the true value
- the estimate is *efficient*, which means that, of all unbiased estimates which are linear functions of the observables, the least squares estimate has the smallest variance. This is called the “Gauss-Markov-Theorem”.

Though these properties are strictly guaranteed only for the linear case, they are still retained in most cases where the function $f_{\tilde{\lambda}}$ can be locally approximated by a linear expansion.

The expression X^2 in equation 1 will follow a χ^2 distribution if the function f_{λ} is (sufficiently) linear and if the measurements m_i follow a normal distribution. This property can be used for statistical tests. In particular the second condition should be always kept in mind, as its relevance will become apparent later.

2.4.2 The Kalman filter technique

The least squares parameter estimation as described in the previous section requires the global availability of all measurements at fitting time. There are cases when this requirement is not convenient, for example in real-time tracking of objects, or in pattern recognition schemes which are based on track following, where it is not clear a-priori if the hit combination under consideration does really belong to an actual track.

The Kalman filter technique was developed to determine the trajectory of the state vector of a dynamical system from a set of measurements taken at different times [23]. In contrast to a global fit, the Kalman filter proceeds progressively from one measurement to the next, improving the knowledge about the trajectory with each new measurement. Tracking of a ballistic object on a radar screen may serve as a technical example. With a traditional global fit, this would require a time consuming complete refit of the trajectory with each added measurement.

Several properties make the Kalman filter technique an ideal instrument for track (and vertex) reconstruction [24, 25, 26]. The *prediction* step, in which an estimate is made for the next measurement from the current knowledge of the state vector, is very useful to discard noise signals and hits from other tracks from the fit. The *filter* step which updates the state vector does not require inversion of a matrix with dimension of the state vector as in a global fit, but only with the dimension of the measurement, leading to a very fast algorithm. Finally, the problem of random perturbations on the trajectory, as multiple scattering or energy loss, can be accounted for in a very efficient way. In its final result, the Kalman filter process is equivalent to a least squares fit.

In this article the implementation and nomenclature from [25, 27] is used, and these documents are referred to for a more detailed explanation of the Kalman filter method. In this notation, the system state vector *at the time k*, i.e. after inclusion of k measurements is denoted by \tilde{x}_k , its covariance matrix by C_k . In

our case \tilde{x}_k contains the parameters of the fitted track, given at the position of the k^{th} hit. The matrix F_k describes the propagation of the track parameters from the $(k-1)^{\text{th}}$ to the k^{th} hit.⁵ For example, in a planar geometry with one-dimensional measurements and straight-line tracks, the propagation takes the form

$$\begin{pmatrix} x \\ t_x \end{pmatrix}_k = \begin{pmatrix} 1 & z_k - z_{k-1} \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ t_x \end{pmatrix}_{k-1} \quad (6)$$

where a subset of the track parametrization in section 2.3.1 has been used. The coordinate measured by the k^{th} hit is denoted by m_k . In general m_k is a vector with the dimension of that specific measurement. For tracking devices measuring only one coordinate, m_k is an ordinary number. The measurement error is described by the covariance matrix V_k . The relation between the track parameters \tilde{x}_k and the *predicted* measurement is described by the projection matrix H_k . In the example in section 2.2.3, the measured coordinate in the stereo view u is

$$H \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} \cos \alpha_{st} & -\sin \alpha_{st} \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} \quad (7)$$

with α_{st} as the stereo angle (45° in the example).

In each filter step, the state vector and its covariance matrix are propagated to the location or *time* of the next measurement with the *prediction equations*:

$$\tilde{x}_k^{k-1} = F_k \tilde{x}_{k-1} \quad C_k^{k-1} = F_k C_{k-1} F_k^T + Q_k \quad (8)$$

and the estimated residual becomes

$$r_k^{k-1} = m_k - H_k \tilde{x}_k^{k-1} \quad R_k^{k-1} = V_k + H_k C_k^{k-1} H_k^T \quad (9)$$

Here Q_k denotes the additional error introduced by *process noise*, i.e. random perturbations of the particle trajectory, for example multiple scattering. We will see later (sec. 5.2) how this treatment works in detail. The updating of the system state vector with the k^{th} measurement is performed with the *filter equations*:

$$K_k = C_k^{k-1} H_k^T (V_k + H_k C_k^{k-1} H_k^T)^{-1} \quad (10)$$

$$\tilde{x}_k = \tilde{x}_k^{k-1} + K_k (m_k - H_k \tilde{x}_k^{k-1})$$

$$C_k = (1 - K_k H_k) C_k^{k-1}$$

with the filtered residuals

$$r_k = (1 - H_k K_k) r_k^{k-1} \quad R_k = (1 - H_k K_k) V_k \quad (11)$$

⁵We assume at this stage a linear system, so that F_k and H_k are matrices in the proper sense. For treatment of the non-linear case see below.

K_k is sometimes called the *gain matrix*. The χ^2 contribution of the filtered point is then given by

$$\chi_{k,F}^2 = r_k^T R_k^{-1} r_k \quad (12)$$

The system state vector at the last filtered point contains always the full information from all points. If one needs the full state vector at every point of the trajectory, the new information has to be passed upstream with the *smoother equations*:

$$\begin{aligned} A_k &= C_k F_{k+1}^T (C_{k+1}^k)^{-1} \\ \tilde{x}_k^n &= \tilde{x}_k + A_k (\tilde{x}_{k+1}^n - \tilde{x}_{k+1}^k) \\ C_k^n &= C_k + A_k (C_{k+1}^n - C_{k+1}^k) A_k^T \\ r_k^n &= m_k - H_k \tilde{x}_k^n \\ R_k^n &= R_k - H_k A_k (C_{k+1}^n - C_{k+1}^k) A_k^T H_k^T \end{aligned} \quad (13)$$

Thus, smoothing is also a recursive operation which proceeds step by step in the direction opposite to that of the filter. The quantities used in each step have been calculated in the preceding filter process. If process noise is taken into account, e.g. to model multiple scattering, the smoothed trajectory may in general contain small kinks and thus reproduce more closely the real path of the particle.

In the equations above, F and H are just ordinary matrices if both transport and projection in measurement space are linear operations. In case of non-linear systems, they have to be replaced by the corresponding functions and their derivatives:

$$F_k \tilde{x}_k \rightarrow f_k(\tilde{x}_k) \quad H_k \tilde{x}_k \rightarrow h_k(\tilde{x}_k) \quad (14)$$

using for covariance matrix transformations

$$F_k \rightarrow \frac{\partial f_k}{\partial \tilde{x}_k} \quad H_k \rightarrow \frac{\partial h_k}{\partial \tilde{x}_k} \quad (15)$$

The dependence of f_k and h_k on the state vector estimate will in general require iteration until the trajectory converges such that all derivatives are calculated at their proper positions. We will continue to call $\partial f_k / \partial \tilde{x}_k$ the transport matrix and $\partial h_k / \partial \tilde{x}_k$ the projection matrix of our system.

The Kalman filter has also been found to be particularly suited for implementation in object-oriented programming language [28].

2.5 Evaluation of Performance

When it comes to quantifying the performance of methods in track pattern recognition, actual numbers will in general strongly depend of the definition of criteria, which comparisons should take into account.

2.5.1 The reference set

Assessment of track finding efficiency requires firstly a definition of a *reference set* of tracks that an ideally performing algorithm should find. Normally tracks will be provided by a Monte Carlo simulation, and the selection of *reference tracks* will depend on the physics motivation of the experiment. Low momentum particles arising from secondary interactions in the material are normally not within the physics scope but merely an obstacle and should be excluded. Particles travelling outside of the geometrical acceptance, for example within the beam hole of a collider experiment cannot be traced by the detector and should be disregarded as well. Also particles straddling the border of a detector and e.g. traversing only a small number of tracking layers will often be regarded as outside of the design tracking volume. A typical convention may be to regard particles which traverse $\mathcal{O}(80\%)$ of the nominal tracking layers as constituents of the reference set.

The definition of the reference set can then be regarded as a definition of effective geometrical acceptance

$$\epsilon_{geo} = \frac{N_{ref}}{N_{total}} \quad (16)$$

with N denoting the number of particles of interest in the reference set and in total.

2.5.2 Track finding efficiency

Definition of the track finding efficiency requires a criterion which specifies whether a certain particle has been found by the algorithm or not. There are two rather different concepts:

Hit matching This method analyzes the simulated origin of each hit in the reconstructed track using the *Monte Carlo truth* information. If the qualified majority of hits, for example at least 70% originates from the same true particle, the track is said to *reconstruct* this particle. This method is stable in the limit of very high track densities, but it requires the Monte Carlo truth information to be mapped meticulously through the whole simulation.

Parameter matching The reconstructed parameters of a track are compared with those of all true particles. If the parameter sets agree within certain limits (which should be motivated by the physics goals of the experiment), the corresponding track is said to reconstruct this particle. This method requires less functionality from the simulation chain, but it bears the danger of accepting random coincidences between true particles and artifacts from the pattern recognition algorithm. In extreme cases, this can lead to the paradox impression that the track finding efficiency *improves* with increasing hit density.

The reconstruction efficiency is then defined as

$$\epsilon_{reco} = \frac{N_{ref}^{reco}}{N_{ref}} \quad (17)$$

where N_{ref}^{reco} is the number of reference particles that are reconstructed by at least one track. It should be noted that this definition is such that a value of one cannot be exceeded, and multiple reconstructions of the same track will not increase the track finding efficiency. One should also control the abundance of non-reference tracks which are reconstructed ($N_{non-ref}^{reco}$): normally the relation

$$\frac{N_{non-ref}^{reco}}{N_{total} - N_{ref}} \ll \epsilon_{reco} \quad (18)$$

should hold, otherwise the reference criteria might be too strict.

2.5.3 Ghosts

Tracks produced by the pattern recognition algorithm that do not reconstruct any true particle within or without the reference set are called *ghosts*. A ghost rate can be defined as

$$\epsilon_{ghost} = \frac{N_{ghost}}{N_{ref}} \quad (19)$$

Since the ghost rate may be dominated by a small subset of events with copious hit multiplicity, it is also informative to specify the mean number of ghosts per event.

2.5.4 Clones

The above definitions for efficiency and ghost rate are intentionally insensitive to multiple reconstructions of a particle. Such redundant reconstructions are sometimes called *clones*. For a given particle m with N_m^{reco} tracks reconstructing it, the number of clones is

$$N_m^{clone} = \begin{cases} N_m^{reco} - 1, & \text{if } N_m^{reco} > 0 \\ 0 & , \text{otherwise} \end{cases} \quad (20)$$

and the *clone rate* becomes

$$\epsilon_{clone} = \frac{\sum_m N_m^{clone}}{N_{ref}} \quad (21)$$

In practice, clones can usually be eliminated at the end of the reconstruction chain by means of a *compatibility analysis* [29].

2.5.5 Parameter resolution

The quality of reconstructed particle parameters and error estimates from reconstruction in a subdetector is essential for matching and propagation into another subsystem. For the whole detector, it determines directly the physics performance. The quality of the estimate of a track parameter X_i is reflected in the *parameter residual*

$$R(X_i) = X_i^{rec} - X_i^{true} \quad (22)$$

From the parameter residual distribution, one can then obtain the parameter estimate bias $\langle R(X_i) \rangle$, and the parameter resolution as a measure of its width. The estimate of the parameter covariance matrix can be used to define the *normalized parameter residual*

$$P(X_i) = \frac{X_i^{rec} - X_i^{true}}{\sqrt{C_{ii}}} \quad (23)$$

which is often called the *pull* of this parameter. Ideally, the pull should follow a Gaussian distribution with a mean value of zero and a standard deviation of one.

2.5.6 Interplay

Results for the individual performance estimators may very much depend on the definitions, so it is advisable to always judge several of the above quantities in combination. For example, the track finding efficiency should be always seen together with the ghost rate, since a less strict definition of the criterion if a track reconstructs a particle will lead to a higher track finding efficiency but also to a higher ghost rate. Also the parameter resolution will tell if the reconstruction criterion is correct, because in case of an inadequately generous assignment, the parameter residuals are likely to show an increased width, or tails from improperly recognized tracks. When parameter matching is used, generous definition of the matching criteria will also increase the track finding efficiency, but reveal itself in a high clone rate.

Excessive tightening of the reference set criteria can potentially also ameliorate the visible track finding efficiency, but it will be at the cost of the effective acceptance, since the total yield of particles with a certain physical signature is proportional to the product

$$\epsilon_{total} = \epsilon_{reco} \cdot \epsilon_{geo} \quad (24)$$

always assuming that relation (18) holds.

3 Global Methods of Pattern Recognition

The task of pattern recognition in general can be described by the illustration in fig. 10. The physical properties of the particles that are subject to measurement are described by a set of parameters, as point of origin, track direction or momentum. Each particle can therefore be represented by a point in the *feature space* spanned by these parameters. The signals the particle leave in the electronic detectors are of a different kind, they are measured hit coordinates the nature of which is governed by the type of device. These coordinates are represented in the *pattern space*. While the conversion from *feature* to *pattern* is done by nature, or by sophisticated simulation algorithms in case of modelled events, the reverse procedure is the task of the combined pattern recognition and track fitting process.

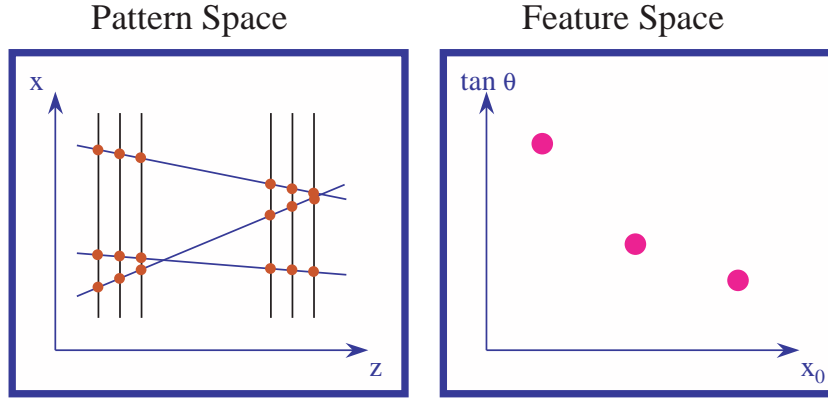


Figure 10: Schematic illustration of Pattern Space (left) and Feature Space (right)

Global methods assess the pattern recognition task by treating all detector hits in a similar way. The result should be independent of the starting point or the order in which hits are processed. This is unlike the *local methods* that will be discussed in section 4, which depend on suitable *seeds* for track candidates. Global methods aim to avoid any kind of seeding bias.

3.1 Template Matching

The simplest method of pattern recognition can be applied if the number of possible patterns is finite and the complexity limited enough to handle them all. In this case, for each possible pattern a template can be defined, for example a set of drift chamber cells through which track candidates in a certain area will pass. Such a technique has been used for the *Little Track Finder*, which was part of the second trigger level of the ARGUS experiment [30], and which worked by comparing the hits in the drift cells of the axial layers to masks stored in

random access memory. This method allowed for basic track finding in a 2D parameter space, the track azimuth and the curvature in the R/ϕ projection, within $20\ \mu\text{s}$. The granularity of the ARGUS drift chamber was moderate, which limited the number of templates that had to be generated. The concept was later extended to the ARGUS *vertex trigger* [31], which used the hits of the micro-vertex detector [32] and generalized the algorithm to three dimensions and four parameters (track curvature being negligible), which allowed to measure the track origin in z to reject background interactions in the beam pipe. This algorithm required the definition of more than 245000 masks, where a five-fold symmetry of the detector had already been exploited.

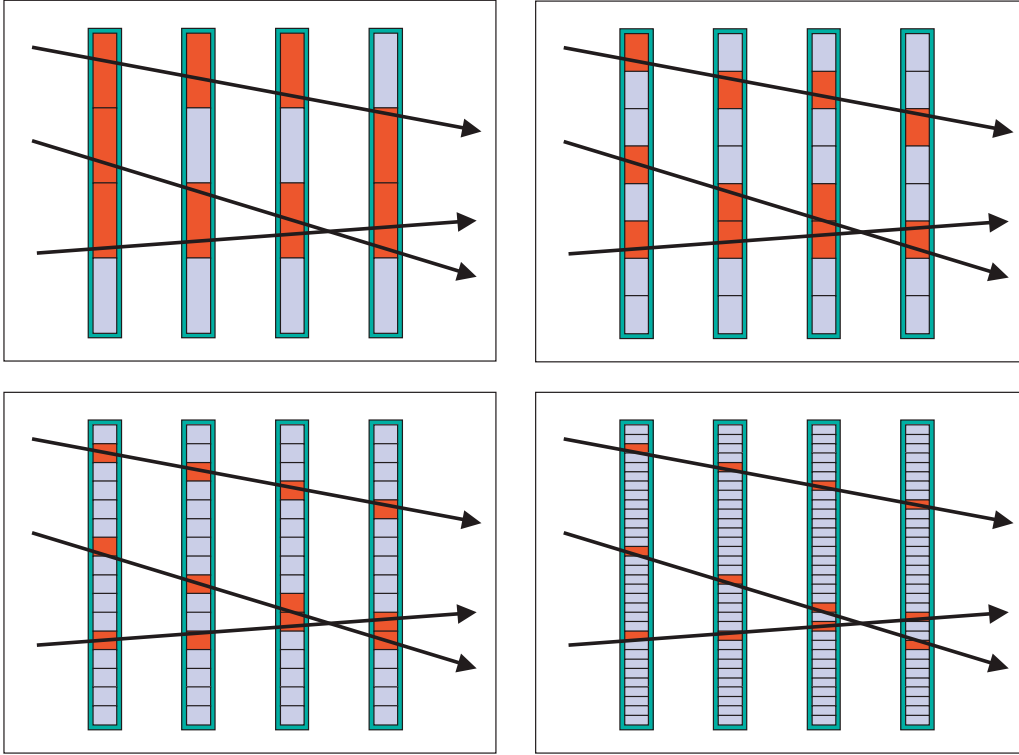


Figure 11: Schematic illustration of the tree-search algorithm: in several steps (in this case four), the track is matched with templates of increasing granularity and resolution. Each step descends into the next level of template hierarchy.

Template matching algorithms are mathematically so simple that they can be hard-wired as *track roads*, provided that the hit efficiency of each element is close to one. Remarkably, the computing time may be independent of the event complexity, since the number of templates to be checked is always the same. However, template matching does not scale very well when the problem requires high dimensionality or granularity. On one side, with increasing granularity the number of templates quickly exceeds limits of feasibility already when storing

them. Also the number of computations increases strongly with a finer resolution of templates. Keeping the granularity low, on the other hand, means that dense situations cannot be resolved, and other methods have to be used to disentangle them.

An elegant solution to both problems is the *tree-search* algorithm, which uses templates of increasing structural resolution that are ordered in a hierarchy [33, 34]. In the first step, the hit structure is viewed at a very coarse resolution with a small set of templates (fig. 11). For those templates that have “fired”, i.e. which match a structure prevalent in the event, a set of daughter templates with finer granularity is applied which are all compatible with the first matched template. This subdivision of templates is iterated until either a matching template on the finest level of granularity is reached – indicating that a good track candidate has been found – or a pattern matched at a certain resolution level cannot be resolved at the next level, in which case it is attributed to a random combination of hits.

The tree-search approach avoids the linear growth of the number of computations with increasing granularity that would develop in a purely sequential search; instead, the computing effort, at least for small occupancy, increases only logarithmically with the number of detector channels. The algorithm becomes even handier when storage of all possible templates can be avoided: in many cases symmetries of the detector can be used to formulate rules how the daughter templates can be derived from the parent at run-time, and how they are connected with the event data. The tree-search algorithm is used for example in the pattern recognition of the HERMES spectrometer, where the final detector resolution of 250 μm is reached in 14 steps [35]. Application of tree-search ideally requires considerable simplicity and symmetry in the detector design, and therefore cannot be easily used in many complex cases. In particular inhomogeneous magnetic fields can complicate the application.

3.2 The Fuzzy Radon Transform

In a very general sense, the observed hit density in the event can be described by a function $\rho(x)$, where x is a very general description of the measured set of hit quantities. In absence of stochastic effects, the expected hit density in the pattern space can be described by an integral

$$\rho(x) = \int_P \rho_p(x) D(p) dp \quad (25)$$

where $D(p)$ describes the prevalent population of the feature space, typically a sum of delta functions centred at the parameters of the particles, and $\rho_p(x)$ is the average response function in pattern space for a particle with parameters p , including all detector layout and resolution effects [36].

Pattern recognition can then be regarded as an inversion of the above integral from a stochastically distorted $\rho(x)$. The Fuzzy Radon transform of the function

$\rho_p(x)$ is defined as

$$\tilde{D}(p) = \int_X \rho(x) \rho_p(x) dx \quad (26)$$

This transformation requires precise knowledge of the response function, in particular the detector resolution. Track candidates are then identified by searching local maxima of the function $\tilde{D}(p)$.

This method shall be illustrated in a simple example with a tracking system consisting of ten equidistant layers in two dimensions without magnetic field. Tracks are parametrized by an impact parameter x_0 and a track slope $t_x = \tan \theta_x$ as defined in sec. 2.3.1. As the measurement is one-dimensional, each hit coordinate gives a linear warp-like constraint in the parameter plane, where the width of the warp reflects the effect of the detector resolution (fig. 12a). For a fictitious situation with three superimposed tracks, the resulting Fuzzy Radon transform is shown in fig. 12b. The three peaks are very pronounced, but development of additional local minima is already visible even in this clean situation.

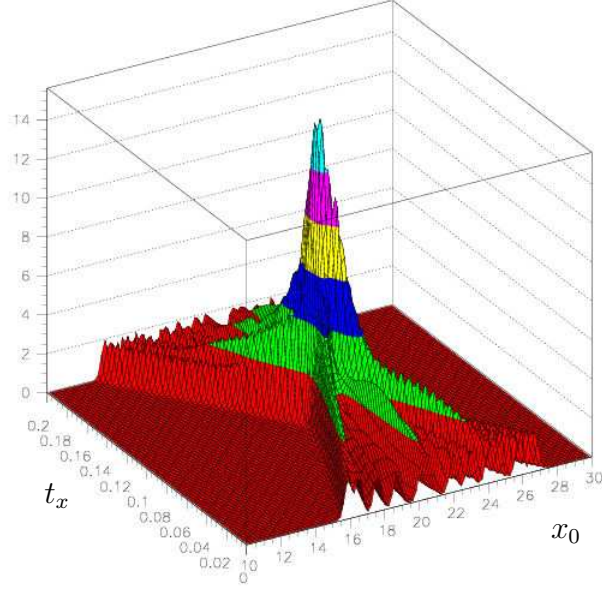
In [36] this method has been explored for a cylindrical geometry in the case of two very close tracks which only differ by a small difference in the curvature value (fig. 13), with additional noise taken into account. Figure 14 shows the resulting Radon transform $\tilde{D}(\kappa, \phi, \gamma)$ as a series of five images around the central values (γ stands for the z speed of the particle which is a measure of the dip angle tangent explained in section 2.3.2), where also the resolution parameter σ has been varied. The images show that the individual tracks can in fact be distinguished (centre image), but it is essential that the assumed resolution parameter matches the real one. It should be noted that automated recognition of the “track signals” in such images would not be a trivial task, and that, for practical purposes, analysis of fuzzy Radon transforms in multi-dimensional parameter spaces are in general very demanding in terms of computing power.

Another generalization of the Radon transform has been investigated in [37].

3.3 Histogramming

As seen in the previous section, the fuzzy Radon transform allows taking the precise detector resolution into account in an elegant manner. In cases where effects of the resolution can be neglected, the response function $\rho_p(x)$ only needs to describe the trajectory, and takes the shape of a delta-function whose argument vanishes for points on the trajectory. This special form of the Radon transform is often called *Hough* transform [39]. The Hough transform of each point-like hit in two dimensions becomes a line; in more generality it defines a surface in the feature space. Completion of the pattern recognition task is thus converted into finding those points in feature space where many of such lines or surfaces intersect, or at least approach each other closely in shape of knots [39].

(a)



(b)

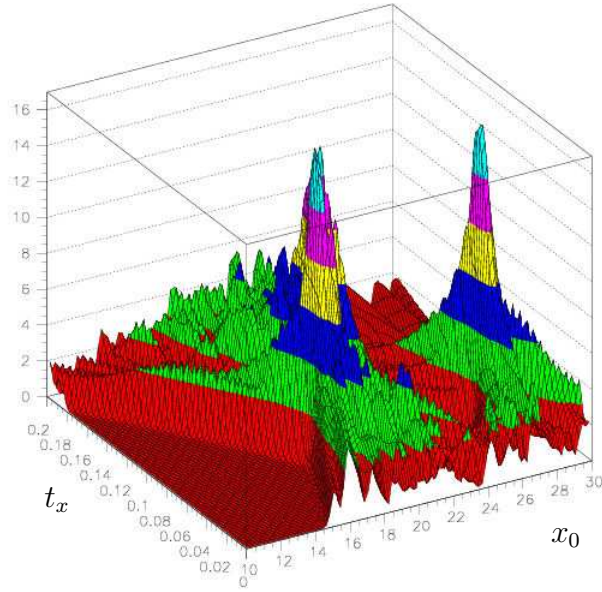


Figure 12: Fuzzy Radon transform $\tilde{D}(x_0, t_x)$ of the hit signals of a single track (a), and in a scenario with three tracks (b), where x_0 and t_x are the track offset and slope.

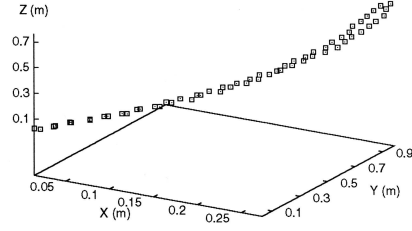


Figure 13: Two simulated tracks differing only by curvature (taken from [36])

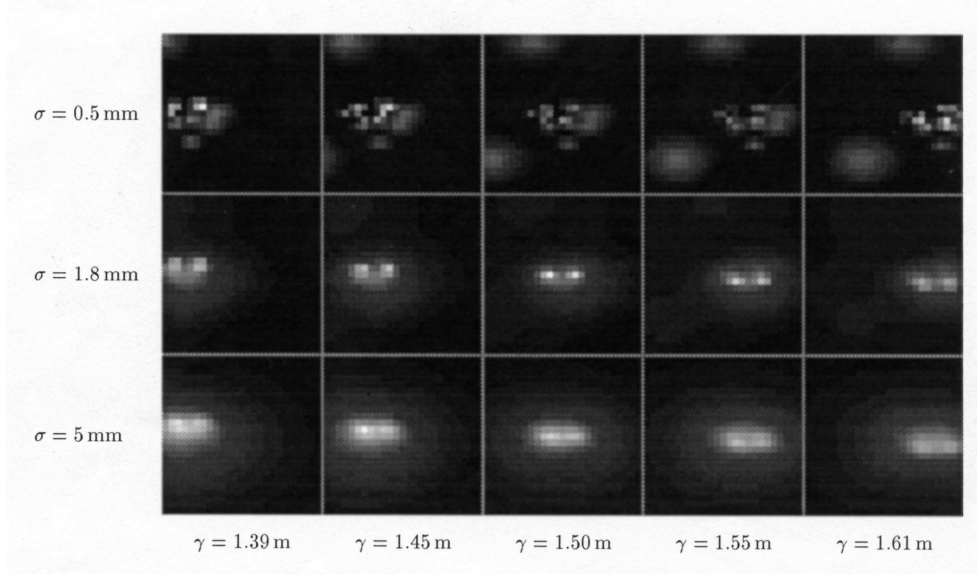


Figure 14: Fuzzy Radon transform of the two tracks in fig. 13 displayed in (κ, ϕ) space, with the third track parameter γ as described in the text (taken from [36]). The transform is shown for three values of the resolution parameter σ in $\rho_p(x)$, where the value in the middle row corresponds to the simulated resolution.

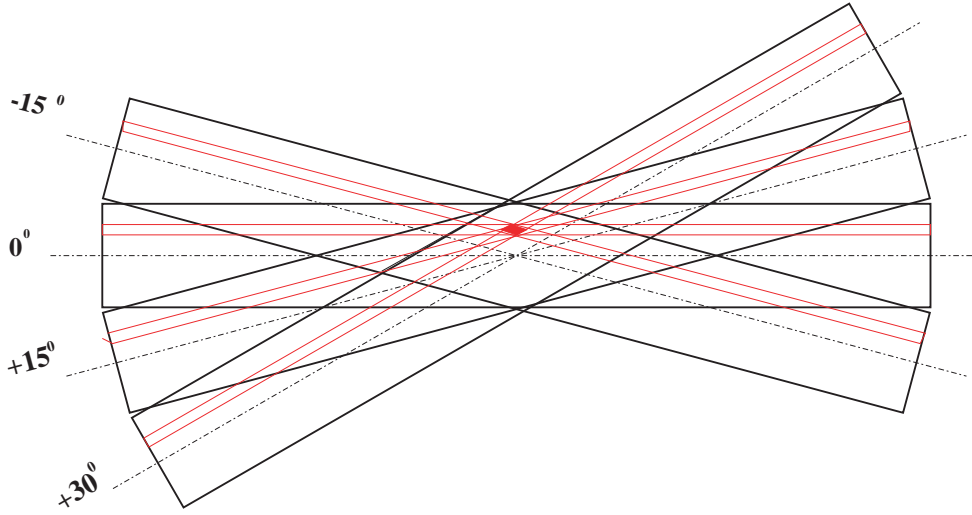


Figure 15: Illustration of wire orientations in the ZEUS straw-tube tracker. In this representation, the beam is oriented vertical to the page, displaced towards the bottom of the page (from [38]).

Histogramming can be regarded as a discrete implementation of the Hough transform. Hit information is converted to a constraint in a binned feature space, and the frequency of entries in a bin above a certain limit is indicative for a track candidate. However, in most tracking devices a single measurement is not sufficient to constrain all track parameters. One solution is then to convert each measurement into a discretized curve or surface in parameter space, and to sample the contribution of all hits in corresponding accumulator cells. An example for such an implementation is shown for the straw-tube tracker (STT) of the ZEUS experiment [38]. This detector system is used as a forward tracker and consists of two superlayers with eight layers of straw tubes each. The straws are arranged in the four different stereo views 0° , $\pm 15^\circ$ and 30° , as illustrated in fig. 15. The 0° straws are oriented such that the point of closest approach to the beam line is in the middle of the straw. Taking the beam spot into account and neglecting the curvature of the segment within the confines of the straw tube tracker, each hit provides an arc-like constraint in the parameter space spanned by polar angle θ and the azimuth angle ϕ . This structure is displayed in the histogram from four views for a single track in fig. 16. The hits from the 0° straws give a transform which is symmetric in azimuth, while the yields from the other views are slightly skewed in correspondence to the stereo angle. The parameters of the track are clearly indicated by the intersection of the four constraints. The resulting histogram is already much more complex in a sample with 10 simulated tracks, where combinatorial overlaps occur (fig. 17).

Another popular way of avoiding the underconstrained case is to combine

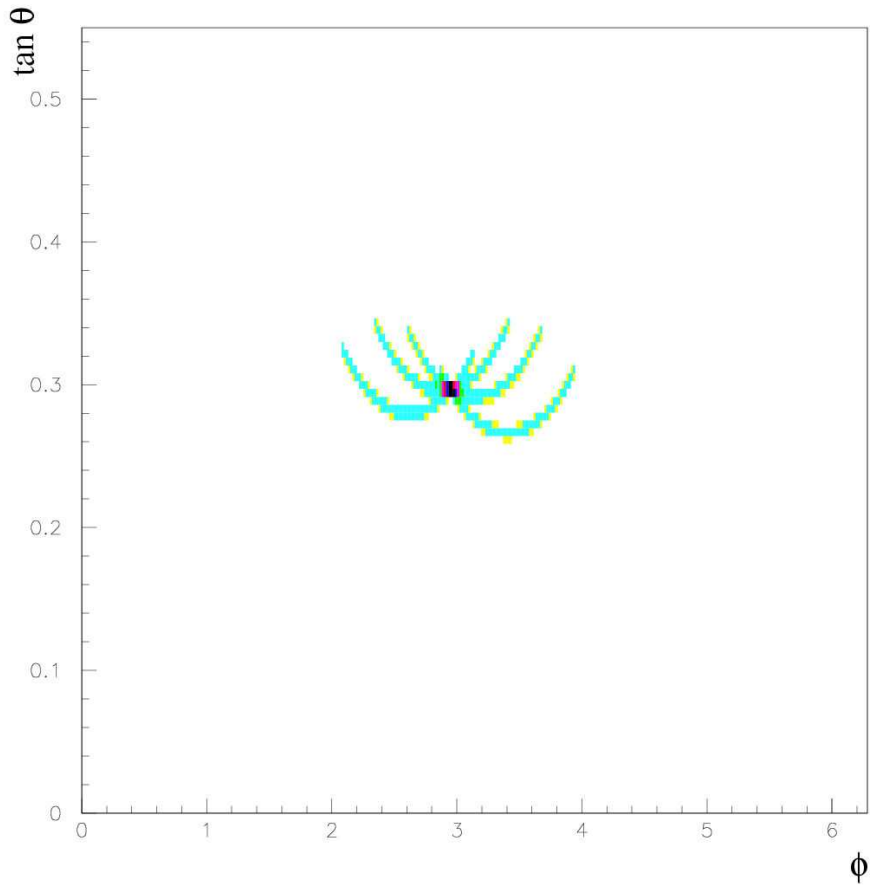


Figure 16: Hough transform of a single simulated track in the ZEUS straw-tube tracker (from [38]).

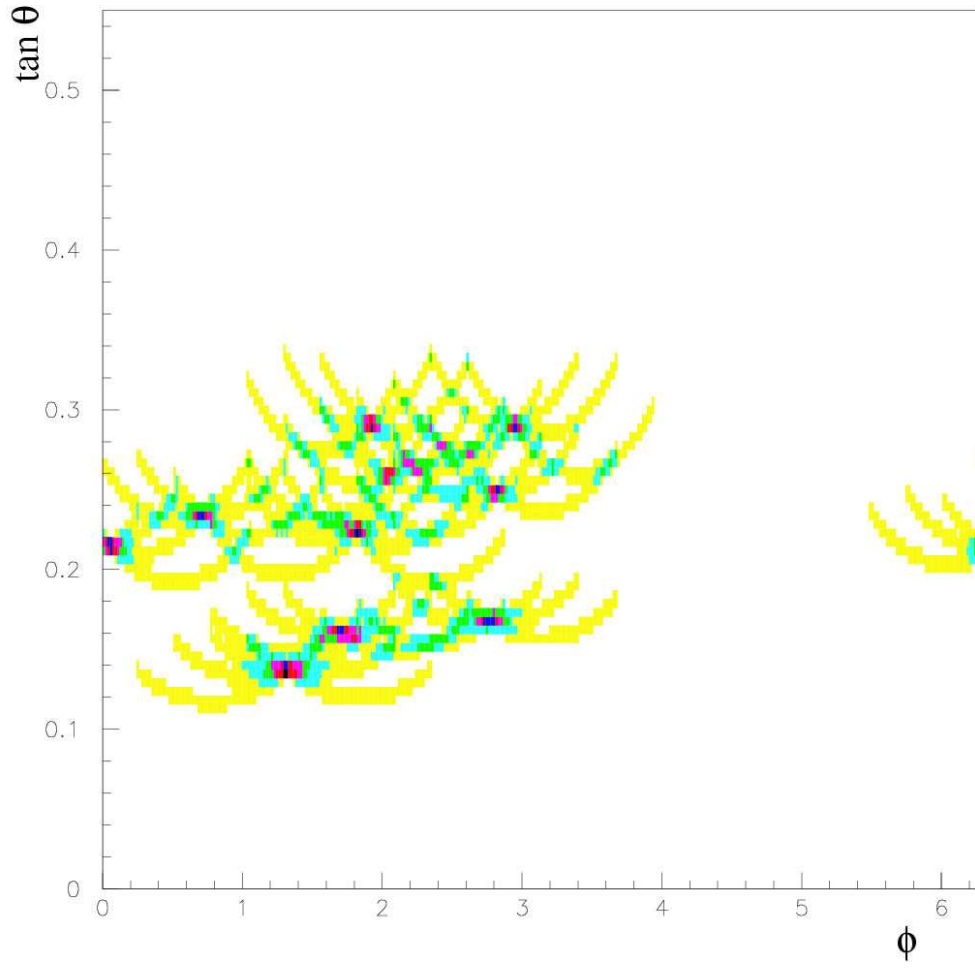


Figure 17: Hough transform of a set of simulated tracks in the ZEUS straw-tube tracker (from [38]).

several hits to track segments before applying the Hough transform. For example, in a 2D *pattern space* without magnetic field, two measured coordinates in the same projection from nearby hits in different detector layers give a straight track segment which represents a point in the *feature space*. Histogramming all segment entries in the feature space should then reveal track candidates as local maxima. This procedure is often referred to as *local Hough transform* [40].

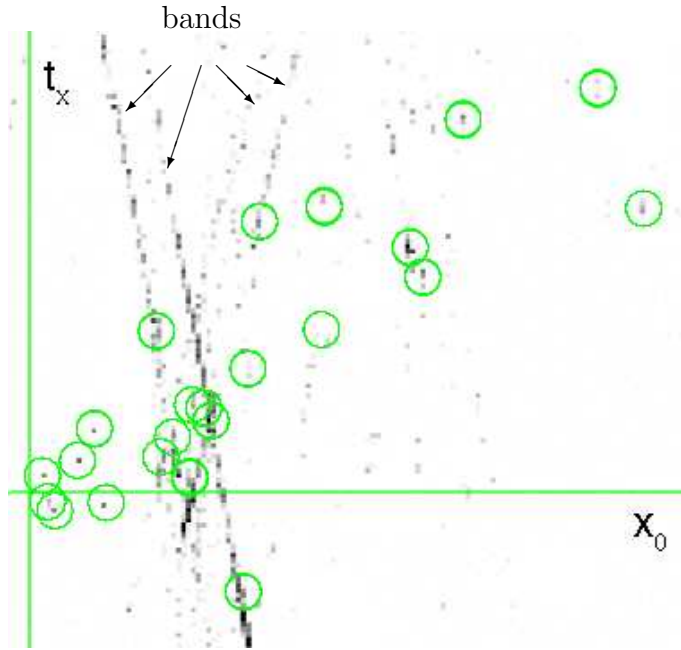


Figure 18: Local Hough transform in a simulated event with five interactions, in the feature space spanned by impact parameter x_0 and track slope $t_x = \tan \theta_x$ (from [41]). The parameters of true particles are illustrated by circles. The colour intensity in each pixel corresponds to the count of segments falling into this square. While the histogram shows the expected enhancements at the *true* parameters of most simulated particles, it also displays artificial structures, indicated as *bands* in the plot that complicate the analysis.

In general, a price has to be paid for this artificial construction of a higher dimension of measurement, since random combinations of hits of different origin lead to ghost segments. The abundance of such contaminations depends strongly on the hit and particle density. A practical example illustrating this problem is shown in fig. 18 (taken from [41]). The geometry corresponds to the “PC” part of the HERA-B spectrometer (see fig. 8), which consists of four tracking superlayers, as indicated in fig. 19a, though in the latter the drawing has been simplified from six to three individual layers per superlayer. A simulated high-multiplicity event with five simultaneous pN interactions has been passed through a local

Hough transform, from which a closeup is shown in fig. 18. The genuine tracks as generated by the Monte Carlo are indicated as circles in the feature space. While enhancements on the histogram are clearly seen at the track parameters of the true particles (indicated by circles), the histogram shows a significant number of bands which are caused by the interference of track patterns. Such interference occurs when several tracks cross the same superlayer of the tracking system within a close distance, as illustrated in fig. 19b for four intersecting tracks: the proximity gives rise to a multitude of combinatorial segments, which have roughly the correct spatial information (x_{SL3}), but a wide range of deviating slopes shadowing the entries with the proper value. These segments enter the histogram with their spatial coordinate transformed to the reference plane relative to which all impact parameters are defined (in this case given by $z = z_{ref}$) in the manner

$$x_0 = x_{SL3} + (z_{ref} - z_{SL3}) \cdot \tan \theta_x \quad (27)$$

The wide spread in the slope $\tan \theta_x$ results in a band in the parameter space, where the tilt of the band

$$\frac{d \tan \theta_x}{dx_0} = \frac{1}{z_{ref} - z_{SL3}} \quad (28)$$

reflects the distance of the superlayer (at z_{SLi}) from the reference plane (at z_{ref}). It is therefore not surprising that in the given detector example with four superlayers, bands of four different slopes can occur.

Even in absence of ghost segments from track overlap, the pattern of track signals in the discretized feature space will in general reflect the underlying layer structure of the tracking system. The local Hough transform is usually based on *short segments*, i.e. those composed of hits in subsequent or at least nearby layers, which has the advantage that the line topology of the track is exploited and the background from random hit combinations is still relatively small. However, due to the small leverage, the angular error can be sizeable, which may impose additional difficulty in identifying the track candidates in the Hough transform. *Long segments* spanning across many layers of the tracking system have the principal advantage of better angular resolution. However, a wide variety of hits have to be combined, so that the number of random combinations increases accordingly. The performance of different approaches has been analyzed in detail in [42]. For the individual application, the optimal choice will depend on the relative importance of resolution and multiple scattering effects.

3.4 Neural Network Techniques

The human brain is particularly skilled in recognizing patterns. It is capable of analyzing patterns in a global manner; it is self-organizing, adaptive and fault-tolerant. It is therefore not surprising that methods have been sought for which

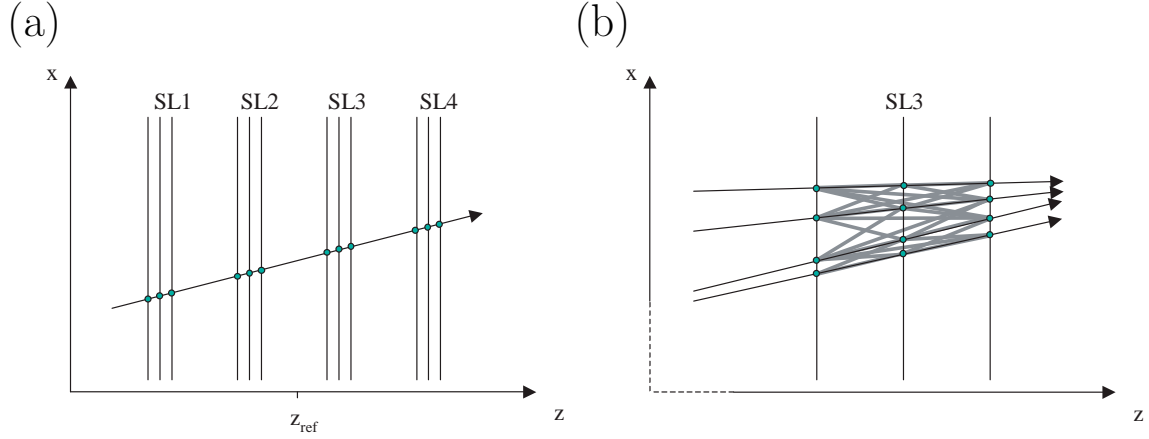


Figure 19: (a) Illustration of the model detector with four tracking superlayers discussed in the text, with the response of a single passing track. (b) Schematic illustration of track segments for a local Hough transform generated from four tracks intersecting in superlayer SL3, showing the abundance of ghost segments compared to the proper ones.

aim at solving pattern recognition problems by means of artificial neural networks. Another intriguing aspect of the human brain is the massively parallel processing of information, which raises hopes that algorithms can be derived which can take full advantage of inherently parallel computing architectures. Because of the wide scope of this subject, this article cannot give a full introduction into this field. A collection of classic papers reprinted is available in [43].

An artificial neuron manifests a simple processing unit, which evaluates a number of input signals and produces an output signal. A neural network consists of many neurons interacting with each other - the output signal of a neuron is fed into the input of many other neurons. While many classification problems can be attacked with simplified layouts, the *feed-forward* networks, track pattern recognition in general uses fully coupled topologies.

3.4.1 The Hopfield neuron

In the Hopfield model [44], each neuron is in general interacting with every other neuron. All interactions are symmetric, and the state of each neuron, expressed by its activation S_i , can only be either *active* (1) or *inactive* (0). The interaction is simulated by updating the state of a neuron according to the activations of all other neurons. The update rule in the Hopfield model sets the new state of a neuron to

$$S_i = \Theta \left(\sum_j (w_{ij} S_j - s_i) \right) \quad (29)$$

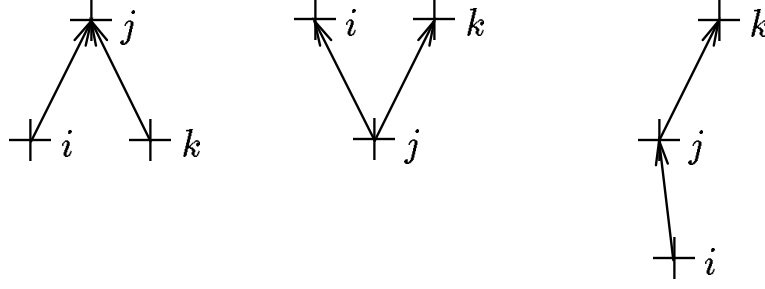


Figure 20: Three typical cases for adjacent track segments in the Denby-Peterson algorithm. The first two combinations correspond to incompatible segments, in the third case, both segments are likely to come from the same track. (from [41])

where the weights w_{ij} determine the strength of each interaction, s_i are threshold values. The theta function $\Theta(\dots)$, whose value is zero for negative arguments and one otherwise, is only the simplest example of an *activation function*, which relates the updated activation to the weighted sum of the other activations. It can be shown [44] that such interactions characterize a system with an energy function

$$E = -\frac{1}{2} \left(\sum_{ij} w_{ij} S_i S_j - 2 \sum_i s_i S_i \right) \quad (30)$$

and that the interaction leads to a final state which corresponds to the minimum of the energy function [44, 45].

3.4.2 The Denby-Peterson method

An adaptation of Hopfield networks to track finding has been developed by Denby [46] and Peterson [47]. The basic idea is to associate each possible connection between two hits with a neuron. Activation of such a neuron means that both hits are part of the same track. It is then essential to define an interaction such that in the global energy minimum only neurons corresponding to valid connections will be active. Interaction is only meaningful with neurons that have one hit in common. An approach to such an energy function is illustrated in fig. 20 [41]: while in the first two cases the neurons (ij) and (jk) represent segments incompatible with the same track and therefore must have a repulsive interaction, the third case is much more track-like and should have an attractive interaction. This desired behaviour can be obtained by an energy function

$$E = -\frac{1}{2} \sum \delta_{jk} \frac{-\cos^m \theta_{ijl}}{d_{ij} + d_{jl}} S_{ij} S_{kl}$$

$$+\frac{1}{2}\alpha\left(\sum_{l\neq j}S_{ij}S_{il}+\sum_{k\neq i}S_{ij}S_{kj}\right)+\frac{1}{2}\delta\left(\sum S_{kl}-N\right)^2 \quad (31)$$

where S_{ij} is the activation of the neuron associated with the segment (ij), i.e. the connection between hits i and j , and θ_{ijl} is the angle between the segments (ij) and (jl). The variables α and δ are Lagrange multipliers preceding terms that suppress unwanted combinations as the first two cases in fig. 20, and fix the number of active segments to the number of hits, N . Track finding is then reduced to finding the global minimum of this multivariate energy function. The interaction is simulated by recalculating the activity of each neuron with the *update rule*, which takes the activations of all other neurons into account.

It is remarkable that the Denby-Peterson method works without actual knowledge of a track model – it favours series of hits that can be connected by a line as straight as possible, but also allows small bending angles from one segment to the next. Thus also curved tracks can be found, provided that a sufficient number of intermediate measurements exists which split the track into a large number of almost collinear segments. The Denby-Peterson algorithm is in particular indifferent about the global shape of the track - a circle and a wavy track with the same local bending angles but alternating directions are of equal value.

One of the first explorations of the Denby-Peterson method has been performed on track coordinates measured by the ALEPH TPC [48]. The algorithm found tracks in hadronic Z^0 decays rather accurately, which may be at least partially attributed to three favourable circumstances: pattern recognition benefits considerably from the 3D nature of the hits measured in the TPC, and equally from the clean event structure and the low occupancy. Moreover, the algorithm is applied such that the initialization activates only neurons that already correspond to plausible connections of hits. The authors of [48] have also investigated the behaviour of the method on events with much higher track numbers, simulated by piling up Monte-Carlo events, and found that the total CPU time of the neural network algorithm is dominated by the initialization of the neurons, which indicates the degree of selection already involved at this stage.

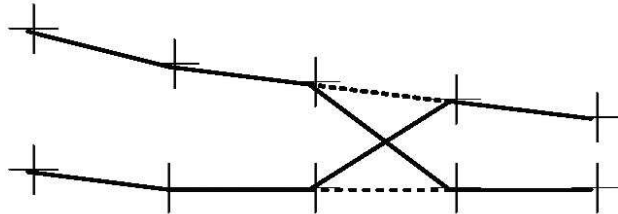


Figure 21: Wrong activations in the case of nearby tracks (from [41]).

The behaviour of the Denby-Peterson method under high track densities has been further investigated in [41] by applying it to a four superlayer geometry resembling the “PC” part of the HERA-B tracker (see fig. 8). These studies found that the classical Denby-Peterson method cannot be relied on to converge safely in cases of nearby parallel tracks. This behaviour is explained in fig. 21: there is no possibility of resolving a cross-wise misassignment, since the system has reached a local energy minimum, and no additional segment can be attached because it would temporarily lead to an illicit branching of the track according to the rules illustrated in fig. 20 and formulated in eq. 31.

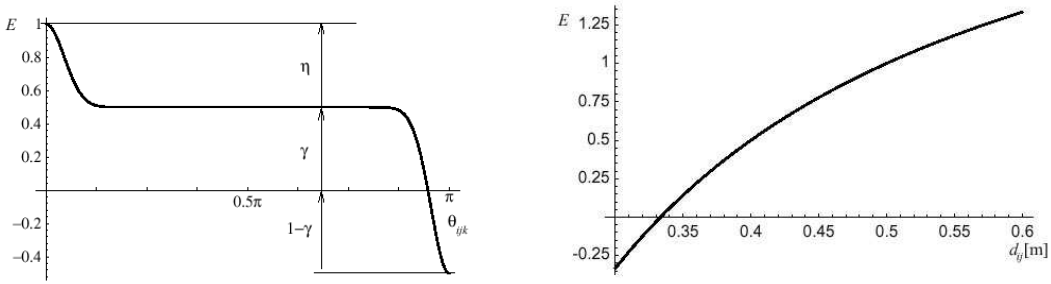


Figure 22: Modified energy function versus angle θ_{ijk} (left) and generalized segment length d_{ij} (right) as used in [41].

The situation can be improved, as shown in [41] by dropping the branching restriction and instead accounting for undesired angles in the cost function, by the replacement

$$\frac{-\cos^m \theta_{ijl}}{d_{ij} + d_{jl}} \rightarrow f(\cos^m \theta_{ij,kl}) \quad (32)$$

where the angle-dependent part is chosen such that only segments with angles close to 180° give a strong negative contribution, and by adding a term proportional to $(\delta - 1/d_{ij})$ for each neuron, which introduces a typical inverse segment length δ into the energy function, where the length of an individual segment d_{ij} is generalized such that the superlayer structure of the tracker is taken into account. (The full definitions are given in [41].) The energy as function of segment angle and length is displayed in fig. 22.

The effect of this variation of the method is visible in fig. 23, which shows the system after one iteration applied to an event with low track multiplicity. At this point, there are still branchings that would not be allowed in the classical Denby-Peterson approach, and which disappear under further iteration. With these modifications the algorithm obtains reasonable efficiency and ghost rate values [49, 41], as displayed in fig. 24.

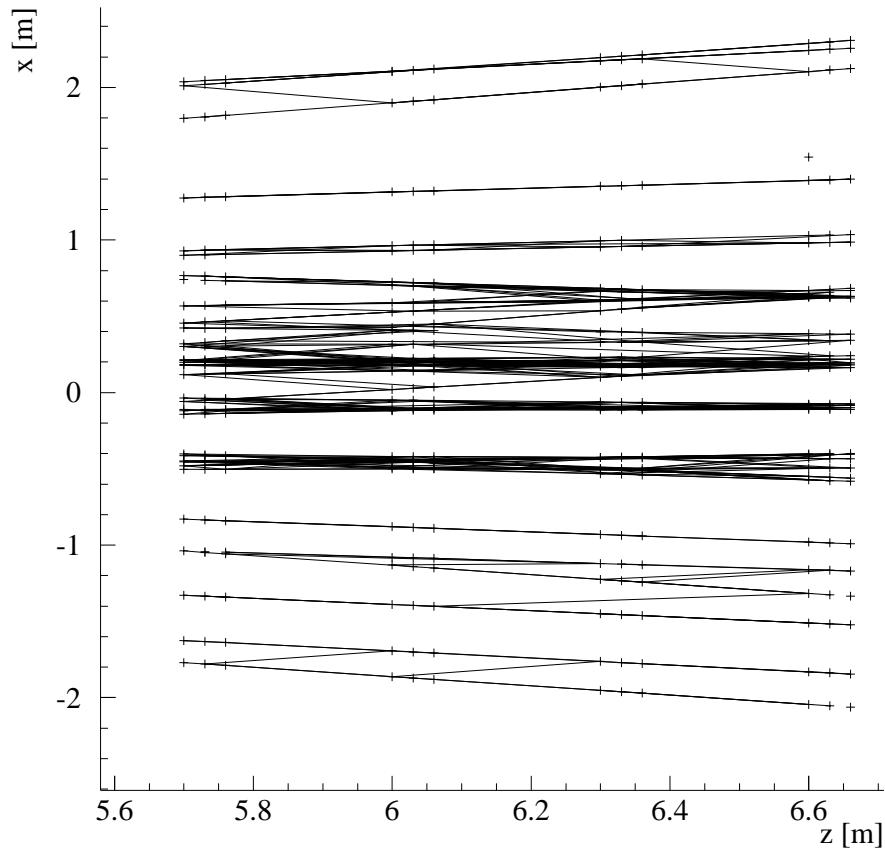


Figure 23: State of the network after one iteration [41]. Crosses denote the locations of the simulated hits.

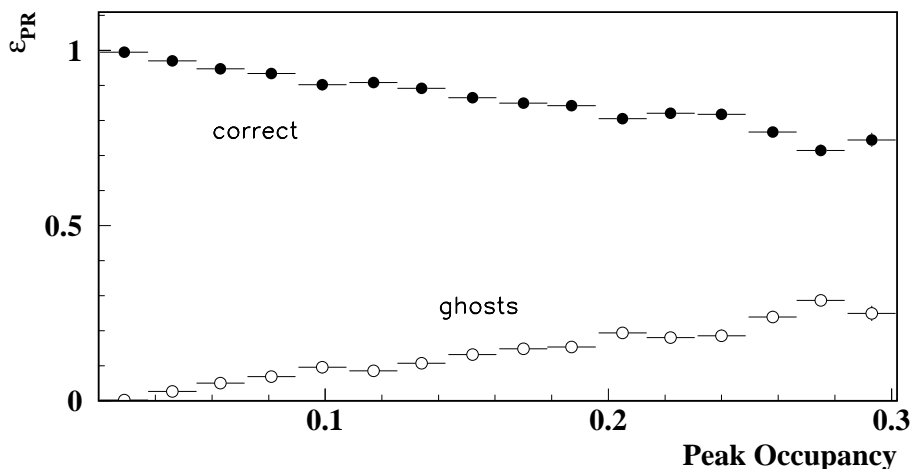


Figure 24: Efficiency and ghost rate for pattern recognition using the modified 2D Denby-Peterson algorithm on simulated events in a fixed target geometry (from [49]).

Several properties of the Denby-Peterson algorithm limit its application at production scale in the general case. The fact that it does not take any explicit track model into account lets it ignore valuable information, which could otherwise help to resolve ambiguous situations. A straight track with random perturbations e.g. is equivalent to a slightly curved track. Neither is there a way to take explicitly the resolution of the detector into account. The computing time per event increases with the third power of the track density, since the number of neurons that have to be generated is proportional to the number of hits squared and the number of non-zero elements in the weight matrix increases with the number of neurons in the vicinity of the track. Perhaps the dominant shortcoming of the Denby-Peterson method is the fact that it does not have a direct extension for finding 3D tracks on the basis of single-coordinate measurements (see 2.2.1), though it is in principle possible to circumvent this problem by first forming space points or segments out of the hits, provided that the ghost combinations are properly eliminated later. Such an approach has been successfully followed in [50], where a method resembling a discrete form of a Denby-Peterson net, referred to as *cellular automaton* [51], was used to select optimal combinations of space points, complemented by a subsequent track following step.

3.4.3 Elastic arms and deformable templates

The above-mentioned limitations of the Denby-Peterson algorithm are overcome with the *elastic arms* algorithm [40, 52], which was introduced by Ohlsson, Peterson and Yuille in 1992. The basic idea can be described as follows: a set of M *deformable templates* is created, which correspond to valid parametrizations of tracks with parameters $\{t_1, \dots, t_M\}$. The number M must be adjusted to the approximate number of tracks in the event. The algorithm should then move and deform these templates such that they fit the pattern given by the positions of N detector hits, which are represented by $\{\xi_1 \dots \xi_N\}$.

As in the Denby-Peterson case, the approach proceeds by formulation of an energy function, whose absolute minimum is at the set of parameters which solve the pattern recognition problem. This requires two elements: an activation-like quantity S_{ia} whose value is one if hit i is assigned to track a , and zero otherwise, and a function $M_{ia}(\xi_i, t_a)$ describing a metric between track template and hit, typically the square of the spatial distance. The energy function can then be defined as

$$\tilde{E}(S, \xi, t) = \sum_{i=1}^N \sum_{a=1}^M S_{ia} M_{ia}(\xi_i, t_a) \quad (33)$$

To avoid trivial solutions, it is necessary to introduce the condition that each hit must be assigned to some template in the form

$$\sum_{a=1}^M S_{ia} = 1 \quad (34)$$

for each hit i . This requirement is called *Potts condition* [53]. One immediate consequence of this condition is the necessity to introduce a special template to which noise hits can be assigned.

The main challenge is then to find the global minimum of the energy function. Since this function tends to be very spiky, as will be illustrated in more detail below, this problem is usually tackled by extending the energy function according to a *stochastic model*, which simulates a thermal motion in the system and smoothens out the spike structure. Search of the minimum starts then at high temperature, and the temperature is successively lowered. At zero temperature, the extended energy function becomes identical to the original one. This technique is called *simulated annealing*.

Instead of the temperature T , normally its inverse $\beta = 1/T$ is used. At finite temperature, the S_{ia} are replaced by their thermal mean values V_{ia} , which take continuous values and lead to a fuzzy hit-to-track assignment. They can be derived from the metric function as

$$V_{ia} = \frac{e^{-\beta M_{ia}}}{e^{-\beta \lambda} + \sum_{b=1}^M e^{-\beta M_{ib}}} \quad (35)$$

where the index b in the sum in the denominator runs over all templates except for the noise template. V_{ia} is called the *Potts factor*. The temperature determines the range of influence for a hit: at zero temperature ($\beta \rightarrow \infty$), the hit is assigned only to the nearest template, with the corresponding V_{ia} equal to one. At higher temperature, the degree of the assignment decreases smoothly with increasing distance. The noise parameter λ represents the symbolic *noise template* which, in the limit of zero temperature, takes over hits that are further than $\sqrt{\lambda}$ away from the nearest genuine template. It is therefore logical to set λ in correspondence to the detector resolution, typically as three or five standard deviations. The term $e^{-\beta\lambda}$ accounts for assignments to the noise template. The Potts factor of the noise template is calculated as

$$V_{i0} = 1 - \sum_{a \neq 0} V_{ia} \quad (36)$$

instead of eq. 35, since the concept of a distance does not make sense here.

The only remaining steps necessary to solve the pattern recognition problem are

1. to find a suitable initialization for the templates, and
2. to find the absolute minimum of the energy function.

It turns out that both are non-trivial in practical applications. Before turning to realistic scenarios, it is very instructive to look at the shape of the energy function in a very trivial example (taken from [41]), which consists of a detector measuring only one spatial coordinate, named x , and a track model consisting only of one parameter for each template. Two hits are considered with coordinates x_1 and x_2 , and two templates with parameters x_a and x_b .

The energy as a function of the template parameters is shown in fig. 25 at a high temperature (the hits being at coordinates $x_a = -1$ and $x_b = +1$). At this temperature, the templates perceive only a blurred image of the hit pattern. The global minimum is at the coordinates in the centre between the hits. When the temperature is lowered to a critical temperature T_c , a saddle point develops (fig. 26), and the previous single minimum splits into two. The critical temperature is related to the coordinates as

$$T_c = \frac{1}{\beta_c} = \left(\frac{x_a - x_b}{2} \right) \quad (37)$$

At very low temperature (fig. 27), two minima have developed at positions corresponding to the two equally valid solutions, $x_a = x_1 \wedge x_b = x_2$ and $x_a = x_2 \wedge x_b = x_1$. The potential ridge at the line $x_a = x_b$ can be interpreted as a repulsive force between the templates [40].

The presence of the noise template parameter λ introduces further local minima into the energy function. An example is shown in fig. 28 with three hits (with

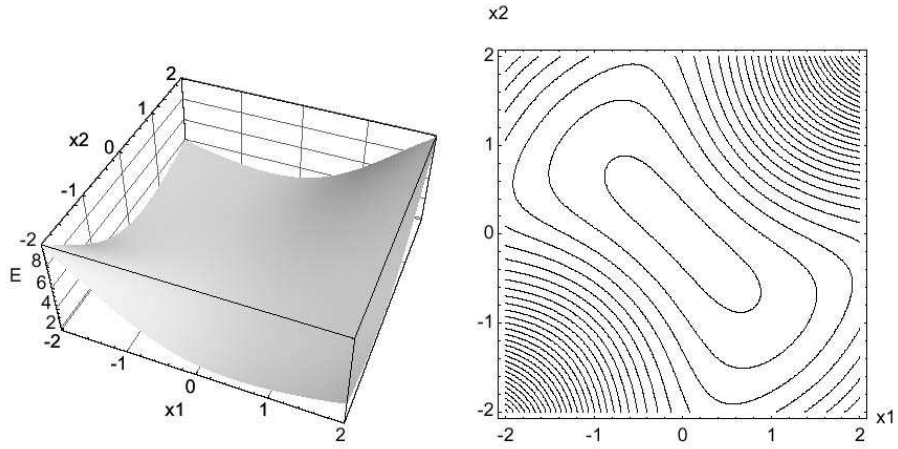


Figure 25: Representations of the energy function of a one-dimensional detector with two hits, as a function of the parameters of two templates x_a and x_b at high temperature [41].

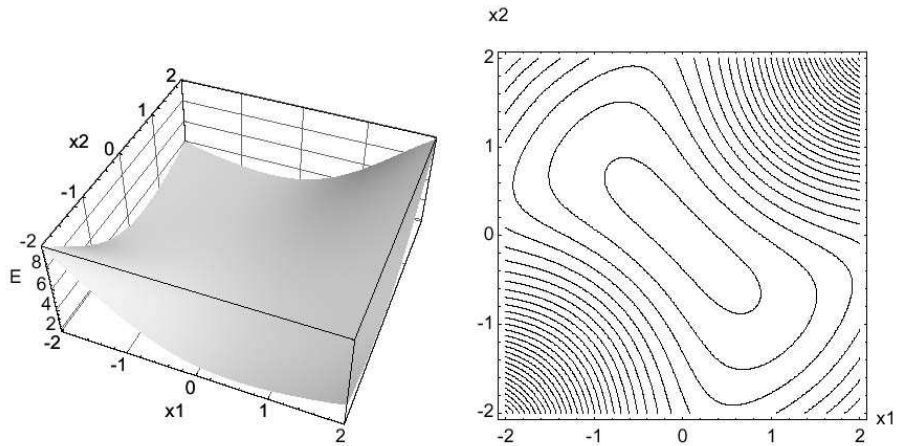


Figure 26: Energy function at critical temperature [41].

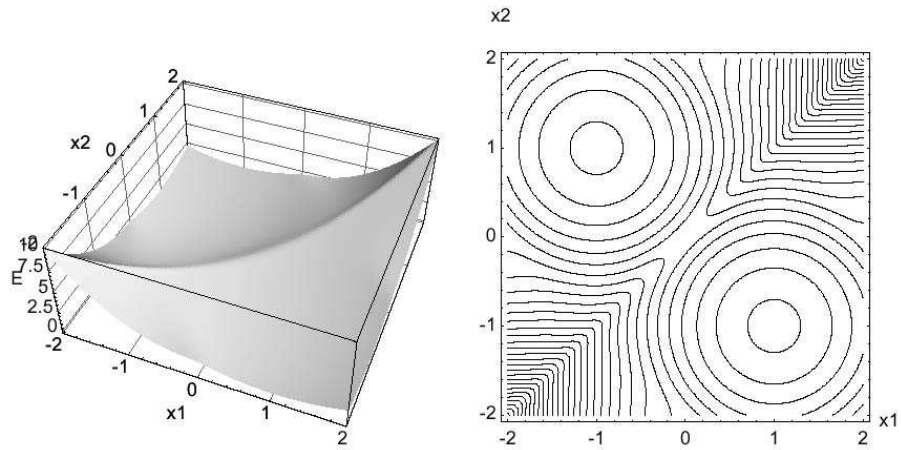


Figure 27: Energy function at low temperature [41].

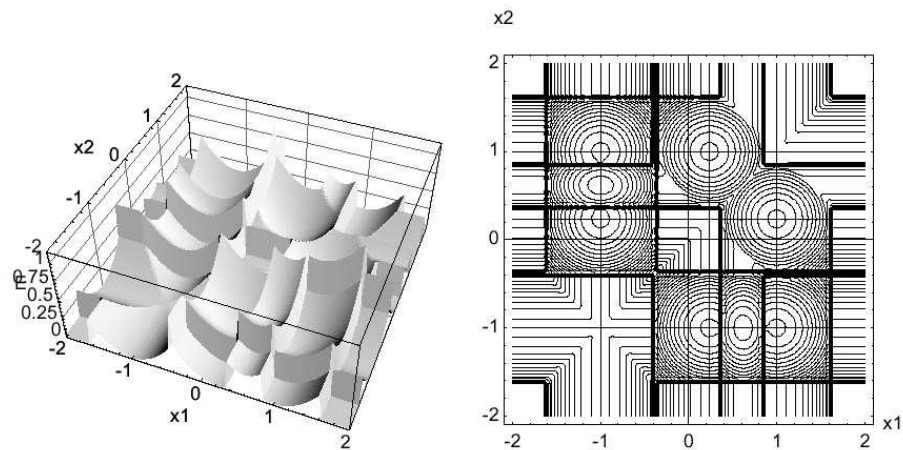


Figure 28: Energy function with three hits at low temperature, with $\lambda = 0.4$ [41]

$x_c = 0.24$) and $\lambda = 0.4$. While the previous solutions are still valid, additional minima appear that correspond to either one or two of the genuine hits being attributed to noise.

The complexity of the energy function for this very simple example is already staggering, and illustrates why initialization and convergence are serious issues.

In their initial study, Ohlsson, Peterson and Yuille [40] applied the method to hits from the DELPHI TPC. Reconstruction was restricted to tracks coming from a vertex spot common to all events, so that track candidates were described by only three parameters, which simplified the situation considerably. The initialization was obtained with a local Hough transform. The moderate hit density allowed performing first the Hough transform in the projection transverse to the magnetic field, searching for track candidates in the space of curvature and azimuth. For each candidate found as a narrow peak in this projection, all hits within a certain neighbourhood were used to calculate the longitudinal tilt angle, which was again histogrammed.

The elastic arms phase then used gradient descent to minimize the energy function at a given temperature. The temperature was lowered by 5% in each step. The Hough transform produced an abundance of templates. The excessive templates were either attracted to noise, or converged to tracks that had already templates associated with them; these had to be weeded out at the end. The result was found to be rather independent of algorithm parameters. The CPU time per event was dominated by the elastic arms step (1 min on a contemporary computer), in contrast to the Hough transform initialization (1 s).

Once more one has to note that pattern recognition in the TPC (here DELPHI's) benefits strongly from the clean event structure with a moderate track density, and the remarkable 3D measurement capabilities of the chamber. An interesting study targeted at much more dense events with 2D measurements has been performed in 1995 [54]. The algorithm was applied to the barrel part of the Transition Radiation Tracker (TRT) of the ATLAS detector, with 40 layers of straw drift-tubes with a diameter of 4 mm and a hit resolution of 150 μm . Since the required hit resolution could only be obtained using the drift-time measurement, the left-right ambiguity had to be resolved. This problem was approached with the elegant method from [55], which introduces energy terms for both left-right assignments (in the nomenclature of eq. 33)

$$\tilde{E}(S, \xi, t) = \sum_{i=1}^N \sum_{a=1}^M S_{ia} \left(s_{ia}^+ M_{ia}^+(\xi_i, t_a) + s_{ia}^- M_{ia}^-(\xi_i, t_a) \right) \quad (38)$$

where the left-right assignment parameters s_{ia}^\pm , which satisfy the condition $s_{ia}^+ + s_{ia}^- = 1$, introduce a repulsive interaction between the alternative left-right assignments, so that a track can only be assigned to one of the two ambiguities of a hit.

The initialization again used a local Hough transform. The minimization phase of the elastic arms step at a given temperature, however, did not rely on

simple gradient descent, but used the *Hessian* matrix, i.e. the second derivative of the energy with respect to the parameters, in a multidimensional generalization of the Newton method. The efficiency was found to be 85% for fast tracks completely contained in the barrel TRT. The efficiency was practically identical to the one of the Hough transform itself, indicating that the elastic arms part did not find any new tracks that had not been properly covered by the initialization. The main advantage of this method is that it validates the tracks found by the Hough transform.

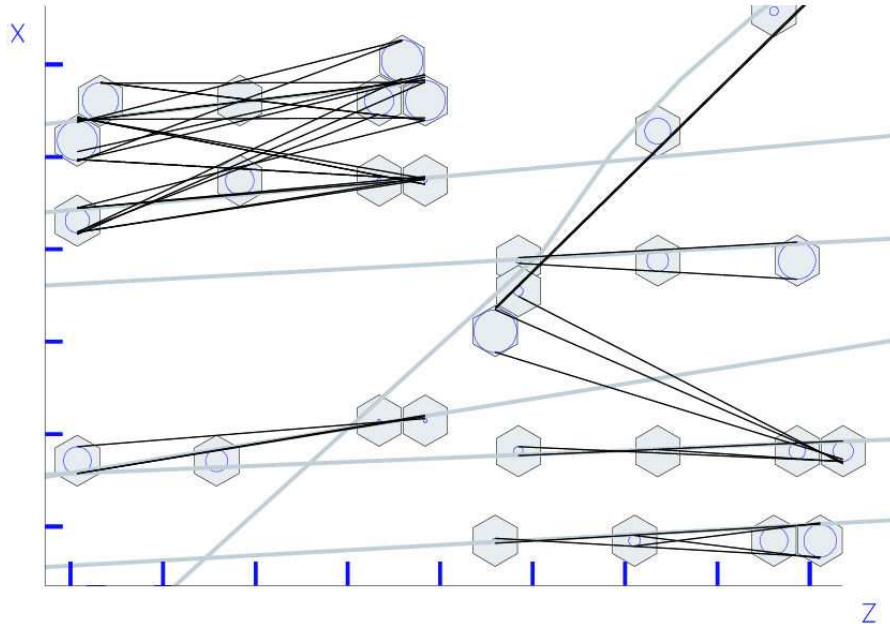


Figure 29: Illustration of segment initialization in the zx projection. The circles are drift distance isochrones of each hit with the drift cell indicated by a surrounding hexagon. The light grey lines are the simulated particles, the black straight lines connecting the hits are the segments produced to initialize the elastic arms algorithm [41].

The track finding capabilities of elastic arms have been further investigated in [41] and [56] with events passed through a full Geant simulation of the “PC” area of the HERA-B spectrometer (see fig. 8). Since the interpretation of the Hough transform turned out problematic in the fixed target geometry under study, a different approach was followed. Track candidates were initialized by searching hit triplets in the 0° projection in each of four superlayers (fig. 29). All triplets with a straight-line-fit yielding $\chi^2 < 3.8$ were accepted, and then matched according to their track parameters. Combinations with triplet segments from all four superlayers were used to initialize the templates in the horizontal plane. The

elastic arms algorithm was then used to perform the pattern recognition together

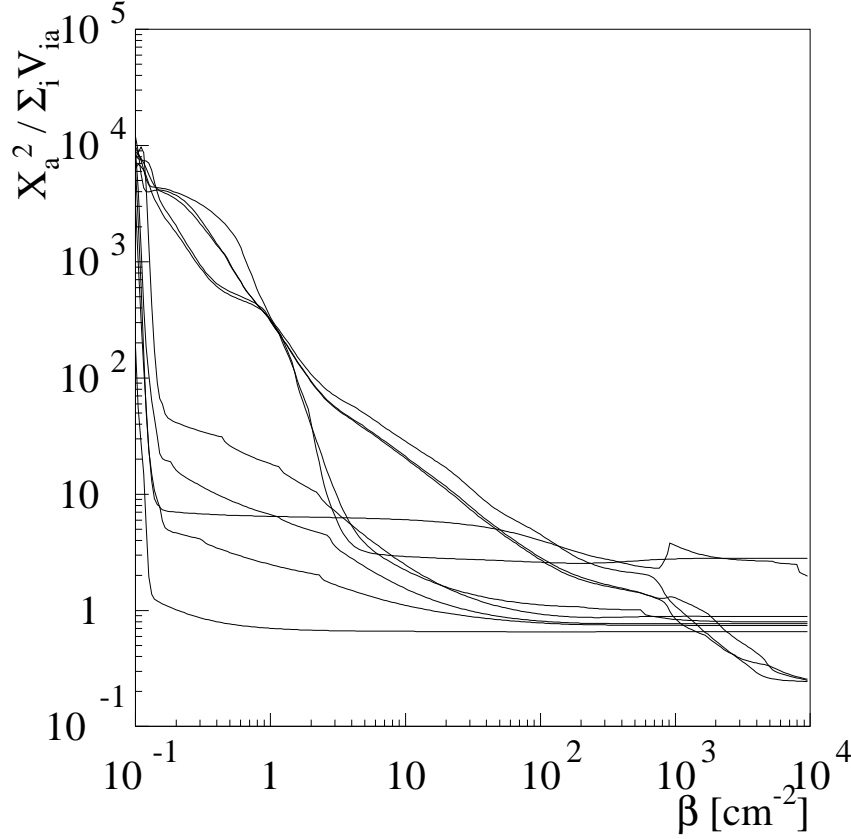


Figure 30: Development of χ^2 with increasing β (corresponding to decreasing temperature) with 10 muon tracks [41].

with the stereo layers, arranging the tracks vertically. Proper operation of this method was shown with test events with ten muon tracks, where the convergence of the tracks in the annealing from a temperature parameter of 0.1 cm^{-2} to 10^4 cm^{-2} is illustrated in fig. 30 by the decrease of the χ^2 per track. While the algorithm was actually performing the task of *vertical pattern recognition* after horizontal initialization, the computing time for the annealing with 10 tracks turned out to be already about 4000s, and it increased at least with the second power of the number of templates. For this reason, dense events with 100 and more track candidates could not be seriously addressed with this method.

For this reason, a subsequent study [56] focused on the reduction of the processing effort. The first major step was the extension of the segment initialisation to 3D. This was achieved by using the segments found from triplets in the xz projection to convert the information from the stereo layers to 3D coordinates: the

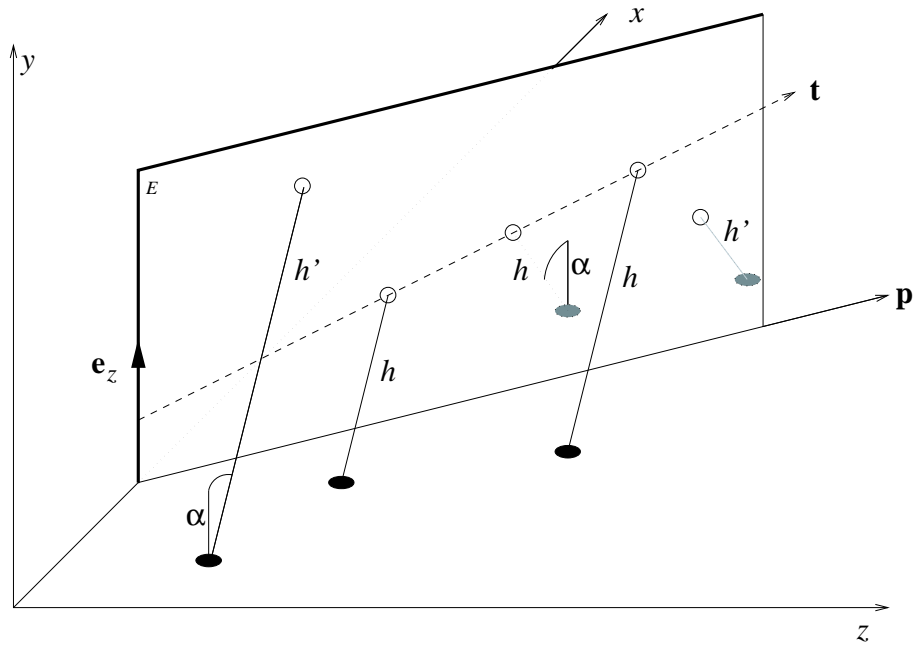


Figure 31: Scheme of converting information from stereo layers to vertical (y) coordinates, by using the horizontal projection track candidate that is indicated by \mathbf{p} . The true track is indicated by a dotted arrow labelled \mathbf{t} . The ovals filled either in black or grey are coordinates measured under the stereo angle $\pm\alpha$ projected into the xz plane. The lines labelled h are stereo hits stemming from the true track, these fall on the trajectory \mathbf{t} when the y coordinate is inferred with eq. 40. Other hits of different origin (labelled h') lead to background hits in the vertical plane (from [56]).

segment in the projection defined a vertical plane in which the track candidate had to be contained (fig. 31). Intersections of stereo wires with this plane lead to indirect measurements in the vertical coordinate y ; the measurement equation

$$u[v] = x \cos \alpha - (\pm)y \sin \alpha \quad (39)$$

was inverted to

$$y = \pm \frac{x \cos \alpha - u[v]}{\sin \alpha} \quad (40)$$

and the triplet and segment finding proceeded with the stereo layers in a similar fashion. The stereo coordinates u and v took drift distance measurements into account, which improved the resolution but lead to left-right ambiguities also in the vertical segment finding.

The second crucial improvement concerned the minimization algorithm within each annealing step. The simplicity of the gradient descent method has made it highly popular for neural network applications, but as already observed in [54], it is by far not the most efficient method. One of its main drawbacks is the fact that its convergence slows down as it approaches the minimum where the surface of the energy function flattens out. On the other hand, large gradients as they can easily occur at lower temperatures (see fig. 28) tend to increase the step size drastically and throw the algorithm completely off the mark. These effects contribute largely to the high computing demands.

It is therefore promising to explore more efficient minimization techniques for high-dimensioned functions [56]. The *Quickprop* algorithm [57] parametrizes the dependence of the energy function on a template parameter $t_a^{(k)}$ (where a is the identifier of the template and k the index of the template parameter) in second order

$$E(t_a^{(k)})_{\{\mathbf{t}_a\}} = c_0 + c_1 t_a^{(k)} + c_2 (t_a^{(k)})^2 \quad (41)$$

and replaces the parameter in iteration step $(i+1)$ with the value at the minimum of the parabola, which is calculated using the gradients of the two previous steps i and $(i-1)$:

$$\Delta t_{a,i+1}^{(k)} = \frac{-\left.\frac{\partial E}{\partial t_a^{(k)}}\right|_i}{\left.\frac{\partial E}{\partial t_a^{(k)}}\right|_i - \left.\frac{\partial E}{\partial t_a^{(k)}}\right|_{i-1}} \Delta t_{a,i}^{(k)} \quad (42)$$

Another more sophisticated minimization method, the *RPROP algorithm* [58], eliminates entirely the dependence of the step width of the gradient by using only its sign. Each component of the template parameter set has its own step width, which is reduced in each step if the sign of the partial derivative has not changed, and somewhat increased if the sign has changed, indicating a step across the minimum.

In application to fully simulated events, the RPROP algorithm turned out to be ten times faster than simple gradient descent. The Quickprop algorithm reduced the computing time by yet another factor of two, but failed to converge properly on about 10% of the tracks, so that the RPROP algorithm was finally chosen for further study [56].

| N_{int} | Segment initialization | | | Elastic arms (incl. initialization) | | |
|-----------|------------------------|-----------|----------|-------------------------------------|-----------|----------|
| | Efficiency | Ghostrate | CPU time | Efficiency | Ghostrate | CPU time |
| 1 | 91% | 38% | 4s | 90% | 3.7% | 15s |
| 2 | 91% | 100% | 14s | 89% | 5.9% | 40s |
| 3 | 89% | 240% | 47s | 87% | 7.5% | 105s |
| 4 | 87% | 440% | 107s | 86% | 10% | 198s |
| 5 | 85% | 1100% | 234s | 83% | 13% | 371s |

Table 1: Efficiency of segment initialization and elastic arms algorithm as compiled from [56], as a function of the number of superimposed interactions, N_{int} . In the elastic arms section of the table, efficiency, ghost rate and CPU time include the effects of the segment initialization.

The segment initialization achieved a track efficiency of 91% for single interactions, which dropped to 85% for five superimposed interactions in an event (tab. 1). The relative efficiency of the subsequent elastic arms phase was always better than 98%, indicating that hardly any of the good tracks the initialization had found were lost. On the other hand, the elastic arms algorithm strongly reduced the rate of ghost tracks prevalent in the initialization. The CPU time consumption, determined on a HP9000/735 processor with 125 MHz clock rate, was still relatively high, but with slightly more than 2 minutes for five simultaneous interactions already in a feasible range. With increasing track density the CPU fraction of the initialization increased steadily and exceeded that of the elastic arms part beyond three superimposed interactions.

The investigations underline that elastic arms can in principle be employed in an efficient manner, but require a very good initialization of the track candidates. This has lead to the general perception that elastic arms should not be used for track finding from scratch, but should rather be seen as a tool to optimize assignment of hits to tracks, to resolve left-right or other ambiguities, or to detect and eliminate outlier hits. A similar philosophy is followed in [50]. A very interesting development in this context is the *deterministic annealing filter* (DAF) [59, 60], which extends the track fit with the Kalman filter with a fuzzy hit assignment and obtains a mathematical equivalent of the elastic arms procedure.

4 Local Methods of Pattern Recognition

While global methods of pattern recognition have the common property to treat all hit information in an equal and unbiased way, simultaneous consideration of all hits can be very inefficient in terms of speed. In fact many detector layouts provide sufficiently continuous measurements so that the sheer proximity of hits makes it already likely that they belong to the same track. This is one of the reasons why *local* methods of track pattern recognition, often called *track following*, are the workhorses of many reconstruction programs in high energy physics.

Track following methods are essentially based on three elements:

- A parametric track model, which connects a particle trajectory with a set of track parameters and provides a method of *transport*, i.e. extrapolation along the trajectory
- A method to generate *track seeds*, i.e. rudimentary initial track candidates formed by just a minimal set of hits which serve as starting point for the track following procedure
- A quality criterion, which allows distinguishing good track candidates from ghosts so that the latter can be discarded

A related variant of track following is the *propagation* of a track candidate found in one part of the tracking system into another, collecting suitable hits on the way. In this case the initial track candidate takes the rôle of the seed.

4.1 Seeds

There are different possible philosophies how seeds can be constructed. This is illustrated in fig. 32, which shows schematically a tracking system with equidistant layers. Starting from the last layer L, where the hit density is lowest, seeds can be obtained by combining the hit with suitable others in the neighbouring layer K (left side). This is the natural choice which exploits the local proximity of hits as a selection criterion. The angular precision of such a short segment is in general limited because of the small leverage, but the rate of fake seeds is relatively small, since most wrong combinations tend to obtain a steep slope that is incompatible with the relevant physical tracks and can be discarded immediately. A completely different alternative is to combine hits for example from the distant layers K and A to construct seeds. These seeds have potentially a much better precision in angle, but the number of choices to be considered is also much higher. The gain of precision can in fact be very limited if the material within the tracker introduces sizable multiple scattering dilution. For the latter reasons, seed combinations from nearby layers are often preferred in practical applications.

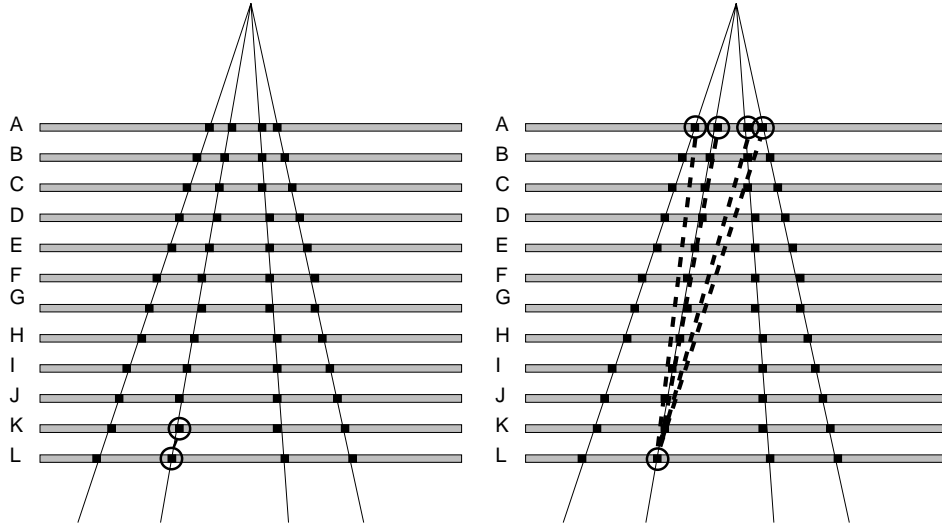


Figure 32: Seeding schemes with nearby (left) and distant layers (right).

Though the number of hits required for a seed is in general dictated by the dimensionality of the parameter space, additional hits can very efficiently decrease the ghost rate of the seeds. Figure 33 shows the construction of seeds consisting of three drift chamber hits each [61]. In this example without magnetic field, only two hits would be minimally needed to define a seed, but the example shows that using hit triplets reduces the combinatorics considerably.

4.2 2D Versus 3D propagation

Many detector layouts allow track following in a projection. For example, drift chambers with many wires parallel and of same length may allow separation of a pattern space that is measured in a plane orthogonal to the wires. This means that parameter propagation during the track following process is far less costly in terms of computations, and that the seeds can be constructed from only two measured hits in the case of a field-free area, or from three hits within a magnetic field. It should be noted that in presence of a magnetic field, a 2D treatment is only possible if the field is oriented parallel to the wire, and homogeneous in wire direction. An example for such an application is the pattern recognition in the ARGUS drift chamber (fig. 34), where the seeds are constructed from three hits in the outer layers, and the track following proceeds towards the beam line [62].

However, pattern recognition in projections cannot avoid that at some point, 3D information must be inferred. This can be achieved by performing track finding independently in all available projections, and then merging compatible projected track candidates into a 3D object. For an unbiased tracking, at least three independent views must exist (see sect. 2.2.3), and each view must possess

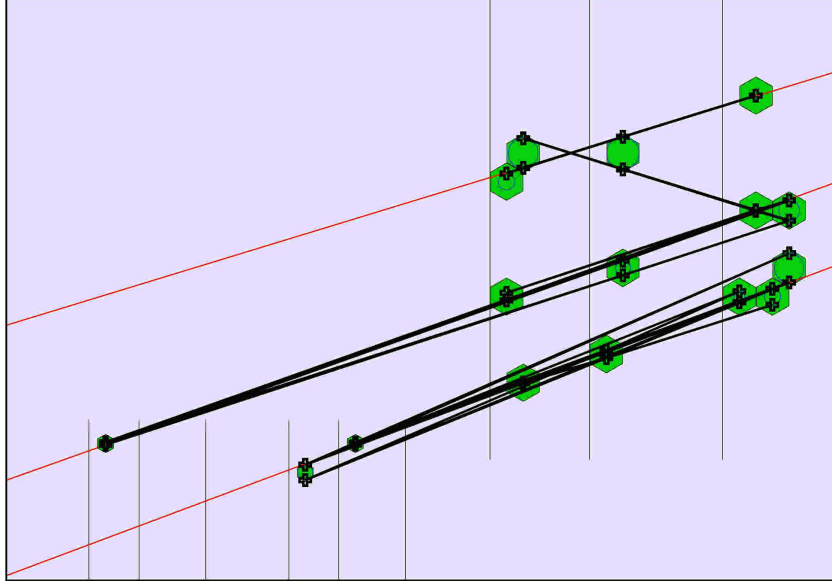


Figure 33: Creating seeds from drift chamber hit triplets. The style of displayed items is similar to fig. 29. Crosses indicate the hit coordinates used to construct the triplets (from [61]).

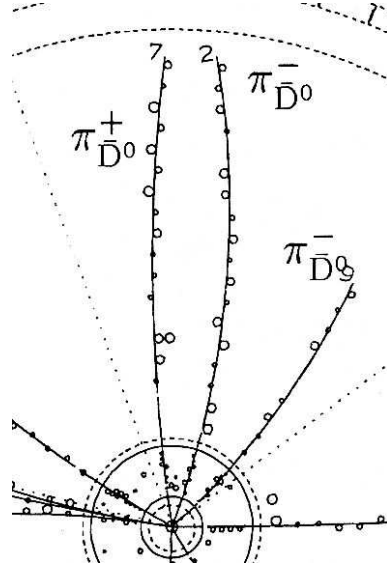


Figure 34: Close-up of the drift chamber area from the ARGUS event display in fig. 1 [2]. The tracks are obtained by a track-following algorithm that proceeds from the outer towards the inner layers.

enough hit information to find the track by itself with good efficiency. A typical symmetric arrangement consists of three views with 0° , 120° and 240° stereo angle, among which all layers are evenly distributed. This approach leads to virtually azimuth-independent track parameter resolutions.

A more economic alternative is a design with an asymmetric layer distribution which is less costly in terms of the number of channels but requires suitable pattern recognition algorithms. It is possible to perform first the pattern recognition in the 0° projection, and then use the resulting track candidate to convert the measurements in the $+\alpha$ and $-\alpha$ layers into the vertical coordinate [62], as already illustrated in a different context in fig. 31. The next step then proceeds with track finding in the vertical projection. In this case, only the 0° projection needs to be equipped with enough layers for a standalone track finding, while the two stereo views are combined and thus the number of layers per stereo view can be smaller. A reasonable scenario for this design comprises 50% of the layers oriented at 0° , 25% in the $+\alpha$ and 25% in the $-\alpha$ projection.

In the case of genuine 3D measurements, 3D seeds can be easily constructed from two hits in the field-free case, and from three hits in the case with magnetic field, which normally will hardly lead to combinatorial problems. This is the situation in the barrel part of the CMS inner detector [63, 64], where three layers of silicon pixel detectors with $150\ \mu\text{m}$ pixel size will be used to initiate track seeds, or in TPCs. In case of intrinsically 2D measurements, 3D seeding has the general disadvantage that the seeds will become rather complex, consisting of 4–5 measurements and under high particle density also many false seeds will be generated. Also left-right ambiguities have a strong impact here: a seed constructed from five drift chamber hits yields 32 ambiguous track parameter sets upon expansion of all possible left-right assignments. Once the seed is constructed, the *track following* step involves many extrapolations of the track parameters which are more costly with the full set of parameters, in particular if the covariance matrix is to be transported as well.

On the other hand, 3D propagation is easier to apply in the sense that the full coordinate information is always available, so that e.g. the decision if the track candidate intersects a particular detector volume or not can be made unambiguously and multiple scattering effects can be accounted for with good precision. The issue of merging the different projections is also avoided.

4.3 Naïve Track Following

The naïve variant shall be discussed here essentially to allow for comparison with the more sophisticated approaches. Starting from a seed, the trajectory is extrapolated to the detector part where the next hit is expected. If a suitable hit is found, it is appended to the track candidate. Where several hits are at disposal, naïve track following selects the one closest to the extrapolated trajectory. This procedure is continued until the end of the tracking area is reached, or no further

suitable hit can be found.

Naïve track following is relatively easy to apply to tracking scenarios with moderate track density and often leads to a reasonable computational effort since the number of hits to be considered is roughly proportional to both the number of layers and the number of tracks. The application to situations with large hit density soon reaches its limitations, since in dense environments, track following runs the risk of losing its trail whenever several possible continuations exist. The main complications can be summarized as follows:

1. Some expected hits may be missing because of limited device efficiency, which will be called a *track fault* in the following. This also includes the case where the hit is existing, but out of expected coordinate bounds, for example because of delta electrons created by the impact of the particle. In drift chambers with single hit readout, the drift time measurement of the followed track can be superseded by another particle passing the same cell closer to the signal wire.
2. Wrong hits may be closer to the presumed trajectory than the proper hits and be picked up in their stead. This can happen easily just after the seeding phase when the precision of the track parameters is still limited, or when some false hits have already been accumulated. A wrong hit may stem from another reconstructable track, from a non-reconstructable low-energy particle, or from detector noise.
3. Left-right ambiguities in wire drift chambers double the number of choices. Especially in small drift cells, e.g. in straw tube trackers, wrong left-right assignments are to some degree unavoidable and need to be coped with.

These aspects can pose a particular problem if the track density is subject to strong variations, e.g. due to a fluctuating number of simultaneous interactions under LHC-like conditions.

4.4 Combinatorial Track Following

This variant is aware of possible ambiguities, and in each track following step, each continuation hit which is possible within a wide tolerance gives rise to a new branch of the procedure, so that in general a whole tree of track candidates emerges. The final selection of the best candidate must be done in a subsequent step, which may involve a full track fit on each candidate. This kind of method is potentially unbeatable in terms of track efficiency, but in general highly resource consuming and therefore only used in special cases with limited combinatorics.

4.5 Use of The Kalman Filter

All track following approaches have to evaluate if a certain hit is compatible with the presumed trajectory and thus suitable to be added to the track candidate. The suitability of a hit should be based on criteria which exploit all the knowledge based on those hits that have been accumulated so far. Not only the track parameters themselves, also their precision needs to be known. The ideal tool in this situation is the progressive fit implemented by the Kalman filter, which has been discussed in section 2.4.2.

The Kalman filter prediction already provides an excellent criterion for hit selection. When a hit is considered to be appended to the track, first the *predicted residual* r_k^{k-1} from equation 9 can be used as a rough criterion. After passing a hit through the filter process (see eq. 10), the *filtered* χ^2 defined in equation 12 is an even more precise measure. In general, the decision power will increase when more and more hits are accumulated in the track candidate. Once the full track is available, the result of the Kalman smoother (eq. 13) can be used to detect and remove further outlier hits.

4.6 Arbitration

In practical applications of track following, means are required to reduce its dependency on the starting point, and to decrease its vulnerability against stochastic influences. This process is called *arbitration*. For example, it is mandatory not to depend on a single option of seeding tracks, which would lead to loss of a track if one of the seeding layer happens to be inefficient, but one will normally use several combinations of layers for seeding. Such redundancy increases the probability to obtain a seed for a track even in presence of device inefficiency. When an expected hit appears to be missing in a layer during propagation, it may be advisable not to discard the candidate immediately, but to proceed further until a fault limit is exceeded. In a case where more than one hit could present a suitable continuation for a track, one might want not to decide immediately for the closest hit but create branches into different candidates which are pursued independently. When a hit appears to be fine for a continuation, the algorithm should account for the possibility that this hit is wrong and the right hit has disappeared for some reason. However, naïvely applied, all these extensions lead to either vast combinatorics, which will explode with increasing hit density, or suffer from ad-hoc limitations. A method to overcome these problems will be detailed in the following.

4.7 An Example for Arbitrated Track Following

This section discusses the *concurrent track evolution* algorithm as an example for an approach to track following with arbitration, which is in detail described in [61, 65].

4.7.1 Algorithm

The basic idea is to allow for concurrency of a certain number of track candidates at any time during the propagation of a certain seed, or even a set of seeds. These tracks are propagated in a synchronized manner from one sensitive tracking volume to the next. At each propagation step for each track candidate, branching into several paths is possible and will in general occur. Multiple branches appear when several continuation hits are consistent with the present knowledge of track parameters, or when more than one tracking volume is within reach. Also the possibility that the expected hit is simply missing, e.g. because of device inefficiency, gives rise to a new branch. Thus the procedure explores the *available paths* for all track candidates *concurrently* which leads to a rapid creation of new track candidates. On the other hand, the number of track candidates should not grow beyond control. This is achieved by applying a quality selection on the whole set of concurrent track candidates after each round of propagation, using suitable estimators for the *quality* of a track. This leads to a favourable timing behaviour even for high multiplicity events. Concurrent track evolution can thus be regarded as a variant of *deferred arbitration* [66]. The actual propagation is based on the Kalman filter.

An illustration of this strategy is shown in fig. 35 taken from [61], which shows a potentially ambiguous situation caused by two nearby tracks T1 and T2 plus a large angle track T3 in five layers of honeycomb drift chambers. For simplicity, it is assumed here that the algorithm discards track candidates with more than one missing hit (*fault*) in a row, and that the maximum number of concurrent candidates is three – in reality, higher limits may be used. It is also assumed that a seed of hits from track T1 has been formed on the right side outside of the figure. The propagation proceeds upstream from right to left. The illustration shows how three parallel candidates arise from different left-right assignments to the two drift chamber hits in layer E, which are propagated through layers D and C – including the tolerance of a fault on track T1 in layer D. In layers B and A, the false paths are discarded because of accumulating too many faults, and the proper reconstruction of track T1 is retained. Track T2 should then be found later with a different seed, while track T3 is likely to be non-reconstructable.

Track following in the naïve sense would always accept the hit with the smallest χ^2 contribution, possibly a good solution when the hit density is small. In the presence of multiple scattering and high hit densities, a wrong hit will frequently have a smaller χ^2 contribution than the proper one, or replace a proper hit which

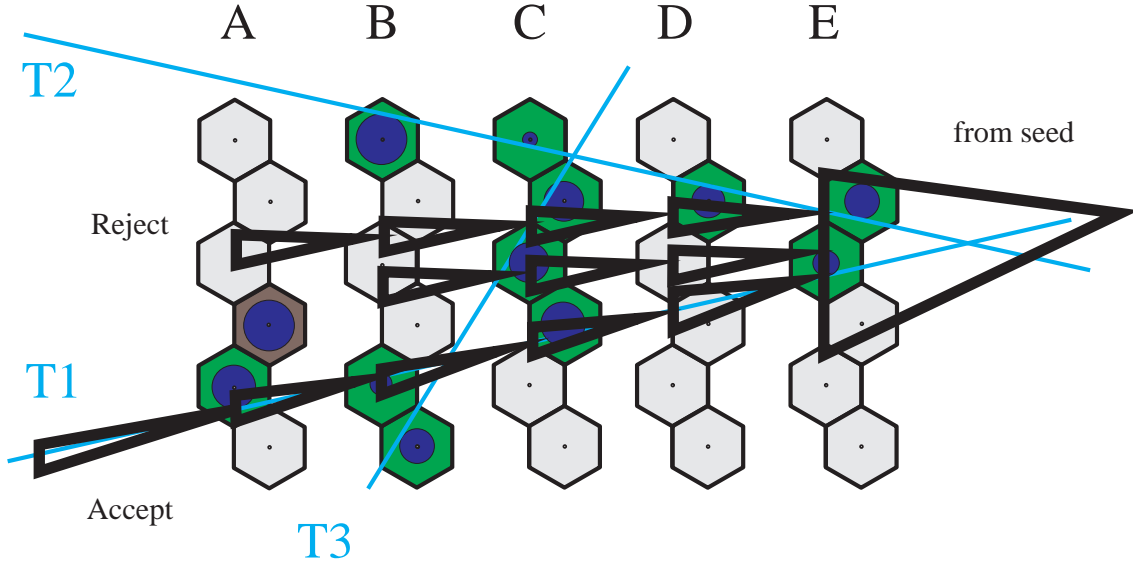


Figure 35: Schematic view of concurrent track evolution in a five-layered part of a tracking system with hexagonal drift cells, which is traversed by three particles, labelled T1, T2 and T3. The simulated drift time isochrones are indicated by circles. The propagation proceeds upstream from the right to the left and starts with a seed of hits from track T1 outside of the picture (from [61]).

is missing due to detector inefficiency, or shadowed by another track passing the same cell. On the other hand, full evaluation of all possible hit combinations would exceed all bounds of computing resources when applied to dense events. Thus, the concurrent track evolution strategy combines the virtues of track following and combinatorial approaches. As will be shown below, the optimization in each evolution step using a quality estimator provides an elegant means to deal with the main problems in high occupancy track propagation.

4.7.2 Parameters

The algorithm is controlled by parameters which determine the selection of hits for propagation of candidates, and for optimization of concurrent candidates on each level. The parameter δ_u^{\max} is the range around the predicted coordinate in the next considered tracking layer, in which continuation hits are searched. The parameter $\delta\chi_{\max}^2$ stands for the maximum tolerable *filtered* χ^2 increment according to eq. 12. Missing hits (faults) are in general tolerated but only a certain number of subsequent faults (N_{Faults}^{\max}) are accepted. The pruning of track candidates after each evolution step is then regulated with absolute and relative cuts. The *quality* of a track candidate can be estimated with a function of the

form

$$Q = f(N_{\text{Steps}}, N_{\text{Faults}}, \chi_i^2, \dots) \quad (43)$$

where N_{Steps} is the number of evolution steps passed so far, and χ_i^2 stands for the contribution of the accumulated hit i to the total χ^2 . If needed, also a bias from the track parameters could be introduced here, which suppresses e.g. tracks that are very steep or have very low momentum. A convenient simple quality estimator is

$$Q = N_{\text{Steps}} - N_{\text{Faults}} - w_{\chi^2} \cdot \sum_i \chi_i^2 \quad (44)$$

which applies a certain malus (in this case 1) for each missing hit, which is equivalent to an ill-matching hit with a χ^2 contribution of $1/w_{\chi^2}$ (in the configuration of tab. 2 equal to 10). Furthermore, cuts are applied relative to the *best* candidate currently in the set: candidates whose quality differs from the best candidate by more than δq_{min} are discarded. Finally, all concurrent track candidates are ranked in decreasing order of quality, and only the first \mathcal{R}_{max} candidates in rank are retained. If propagation cannot be continued though the end of the tracking system is not reached, this may have a natural reason, e.g. the particle may have been stopped or decayed in flight. In such cases, the best remaining track candidate on the last level is kept if it comprises at least a certain minimum number of hits, $N_{\text{Hits}}^{\text{min}}$.

4.8 Track Following And Impact of Detector Design Parameters

The practical behaviour of such an algorithm, as it has been developed for the HERA-B spectrometer has been studied in [61], including an investigation of the impact of detector design and performance on the pattern recognition capability. As the experiment has never routinely taken physics data at the high design interaction rate of 40 MHz, the results have been obtained from a full Geant simulation with on average five superimposed pN interactions, one of them containing beauty hadrons. As seen in fig. 8, the inner part of the HERA-B main spectrometer acceptance within about 25 cm radius from the beam line is covered by micro-strip gaseous chambers (MSGC), while the outer part is instrumented with Honeycomb drift chambers [13, 14, 15]. The *pattern tracker* consists of four superlayers outside of the magnetic field, which consist of 6 individual layers each (the area marked “PC” in fig. 8), except for the inner part of the two middle superlayers that have only four layers each. Half of the layers measure a horizontal coordinate (0° orientation), the other half are arranged at ± 100 mrad stereo angle. The seeds were produced from hit triplets in the hindmost two superlayers for upstream, and in the foremost two superlayers for downstream propagation

| Parameter | Value | Parameter | Value |
|--|-------|--|-------|
| $N_{\text{Hits}}^{\text{min}}(\mathbf{x})$ | 9 | $N_{\text{Faults}}^{\text{max}}(\mathbf{x})$ | 2 |
| $N_{\text{Hits}}^{\text{min}}(\mathbf{y})$ | 9 | $N_{\text{Faults}}^{\text{max}}(\mathbf{y})$ | 2 |
| $\delta\chi_{\text{max}}^2(\mathbf{x})$ | 8 | \mathcal{R}_{max} | 5 |
| $\delta\chi_{\text{max}}^2(\mathbf{y})$ | 16 | w_{χ^2} | 0.1 |
| δq_{min} | -1 | | |

Table 2: Table of parameters used in the implementation in [61]

(fig. 33). Track finding was performed first in the 0° projection, then continued in the combined stereo layers, where the vertical coordinates were determined using the horizontal projection of the track candidate with the method explained in sec. 3.4.3 (see eq. 40 and fig. 31).

The algorithm parameters used are summarized in table 2. The parameters allow for a delicate adjustment of balance between the extremes of naïve track following ($\mathcal{R}_{\text{max}} = 1$), where always the apparently best path is followed, and combinatorial track following ($\mathcal{R}_{\text{max}} = \infty$), which retains all paths. The detailed simulation allowed to study some principal effects of tracking system properties on pattern recognition parameters which will be shown in the following.

4.8.1 Influence of detector efficiency

Figure 36 shows how the hit efficiency of the detector devices affects the pattern recognition performance on tracks emerging from B decays. Above $\epsilon_{\text{HIT}} = 95\%$, the hit inefficiency is well compensated by the algorithm (operating with $N_{\text{Faults}} = 2$), resulting in an excellent track finding performance. Smaller hit efficiency leads to sizeable loss in the fraction of detected particles.

4.8.2 Effect of detector resolution

The influence of the spatial resolution is shown in fig. 37. The simulated resolutions of outer and inner tracking system were varied independently. It is interesting to see that the efficiency degrades only slowly with the resolution being increased up to 1 mm. The slight drop in efficiency at $100 \mu\text{m}$ in fig. 37a is an artifact due to numerical approximations. Both figures indicate that the effect of resolution on track finding efficiency should not be overrated. Much stronger is the effect on the ghost rate, the plots underline that a good resolution helps considerably to suppress fake reconstructions.

4.8.3 Influence of double track separation

The simulation of the inner tracker devices allowed varying of the double track resolution, i.e. the distance down to which nearby tracks can be resolved as

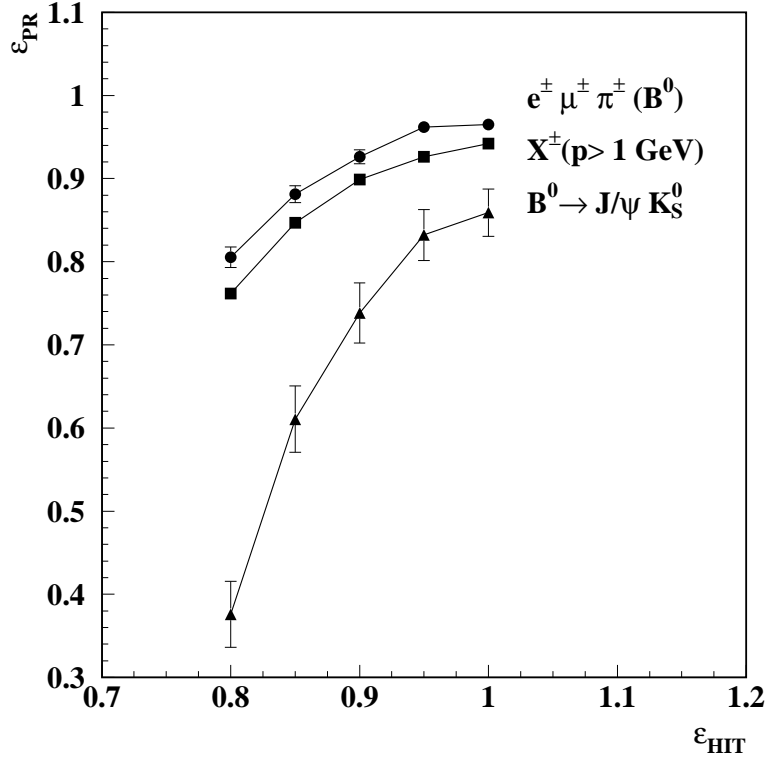


Figure 36: Pattern recognition efficiency for different values of the hit efficiency on simulated events consisting of one pN interaction with a B^0 meson with the decay chain $B^0 \rightarrow J/\psi K_S^0 \rightarrow \ell^+ \ell^- \pi^+ \pi^-$, where $\ell^+ \ell^-$ can be a pair of muons or electrons, superimposed with on average four unbiased inelastic interactions. The filled squares show the track finding efficiency for charged particles with momentum above 1 GeV, the filled circles are for particles from the B decay. The triangles indicate the combined efficiency of all four B decay particles [61].

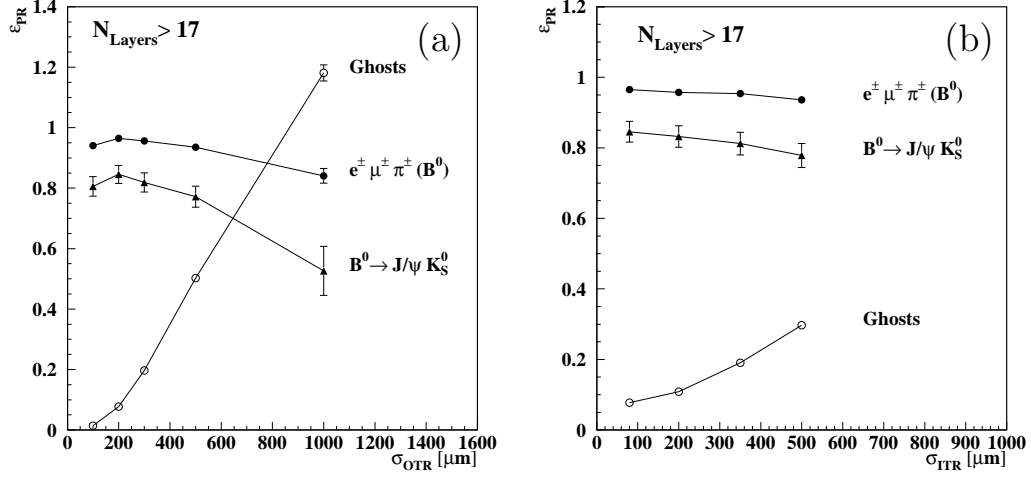


Figure 37: Pattern recognition efficiency for different outer (a) and inner tracker resolutions (b), for particles from the $B^0 \rightarrow J/\psi K_S^0$ decay mode as detailed in the caption of fig. 37. Only tracks passing at least 17 out of 24 possible tracking layers were considered. Also the ghost rate is displayed (open circles) [61].

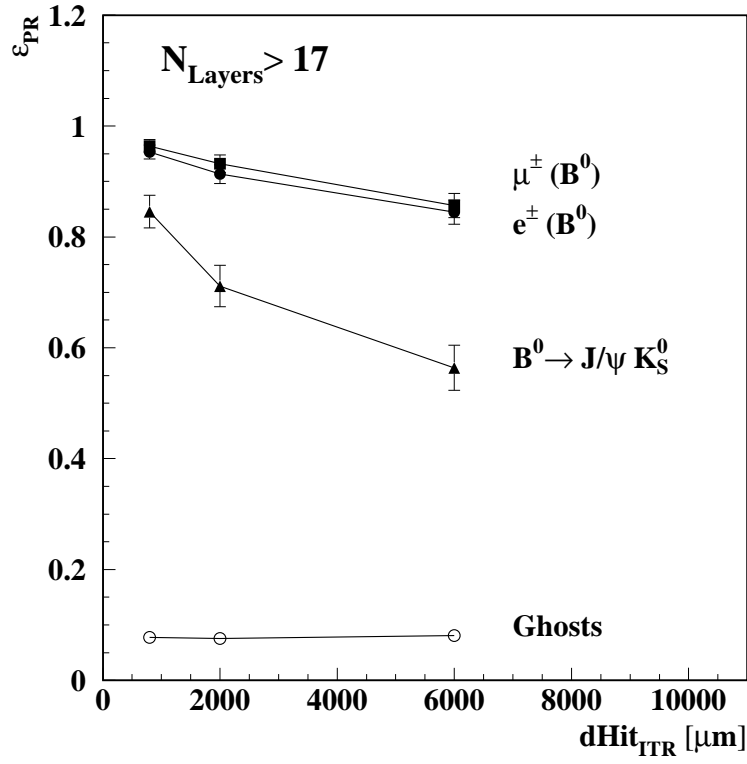


Figure 38: Pattern recognition efficiency for different double hit resolutions of the inner tracker for particles from the $B^0 \rightarrow J/\psi K_S^0$ decay mode as detailed in the caption of fig. 37. Also the ghost rate is shown (open circles). Only tracks passing at least 17 out of 24 possible tracking layers were considered [61].

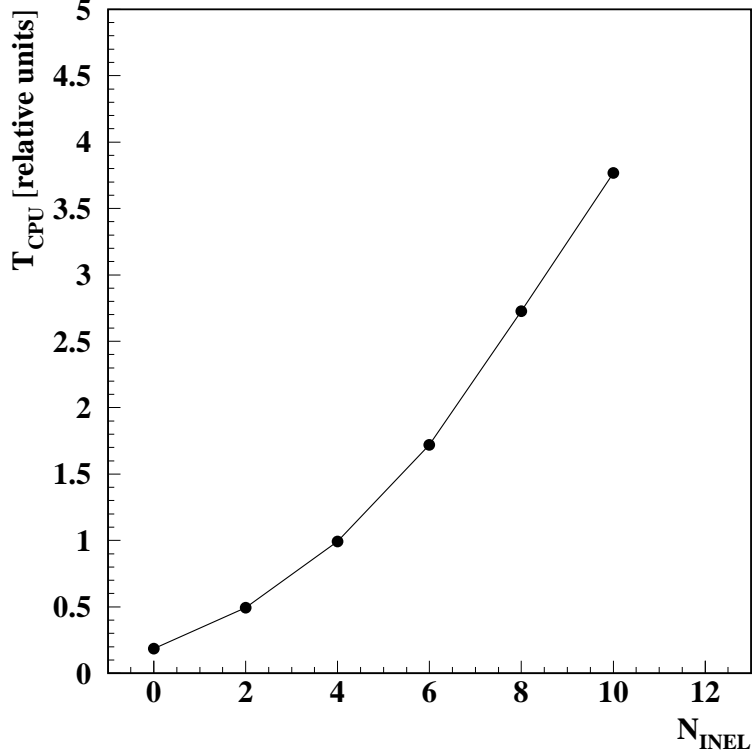


Figure 39: Mean computing time per event as a function of the number of inelastic interactions superimposed to one $b\bar{b}$ interaction [61], normalized to the value at $N_{INEL} = 4$.

individual hits in a device. In micro-strip gaseous chambers (MSGC) as they are used by HERA-B, the double track separation distance is in general larger than the resolution, since it depends on the cluster sizes. As visible in fig. 38, the efficiency drops significantly with double track resolutions worse than $800 \mu\text{m}$.

4.8.4 Execution speed

As already seen in sec. 3.4.3, the CPU time consumption is an essential aspect for the selection of a pattern recognition algorithm. The concurrent track evolution algorithm was tested on the same geometry and event type as the elastic arms algorithm implementation (see tab. 1), and required on average 4s per event with four superimposed inelastic interactions, compared to several minutes for the elastic arms method on the same type of processor. Also the behaviour with increasing track density is important, since steep increases with a sizable power of the track multiplicity, as they may arise from combinatorics, can have a very

negative impact on use of a reconstruction program at production scale. Figure 39 shows the average computing time per event normalized to that for the nominal four superimposed inelastic interactions. At high interaction multiplicity, the computing time per event settled rather gracefully on a roughly linear dependence, indicating a constant amount of time per track, at an acceptable loss of efficiency, which can be considered a good-natured behaviour. With the speed shown above, the algorithm is fast enough to be used in quasi-online reconstruction [67].

4.9 Track Propagation in a Magnetic Field

In general the above track following strategy can be applied also within a magnetic field. The main difference is that the transport function in eq. 8 becomes non-linear, and the transport matrix becomes a local derivative as displayed in eq. 15. If the field is homogeneous, or if inhomogeneity can at least be neglected within typical transport distances, the transport function and matrix can usually be expressed analytically.

In many cases, however, the field is neither homogeneous nor describable in an analytic expression, instead, it is parametrized in terms of a field map, which has been measured with Hall probes, or computed by means of a field simulation program. In this case, numerical methods have to be used to derive the transport function. A very suitable method is the Runge-Kutta procedure [68], which integrates the equations of motion by expanding the trajectory up to a certain order and sampling the field at a series of intermediate points, which are chosen and weighted such that all powers of the errors below a certain order cancel. Even this procedure meets considerable challenges when the field varies strongly and a very high precision, matching the detector resolution, must be warranted. In this situation, an embedded Runge-Kutta method with adaptive step size can help: the next highest order of Runge-Kutta is compared with the preceding one and the difference serves as an error estimate, which is then used to adjust the step size.

Application of the Kalman filter does not only require a transport function for the track parameters, but also the derivative matrix of the new parameters with respect to the old is needed (see eq. 15). Calculation of this derivative matrix can be efficiently performed within the same Runge-Kutta framework that is used for the parameter transport itself [69].

An extension of the concurrent track evolution algorithm for track following in the magnetic field has been developed and tested on the HERA-B geometry in [65]. Track segments found in the field-free part of the spectrometer were followed upstream through the inhomogeneous field of the magnet tracker. Figure 40 shows an event display with simulated tracks including a B decay reconstructed with this method. The algorithm achieved a high track propagation efficiency in spite of the large track density.

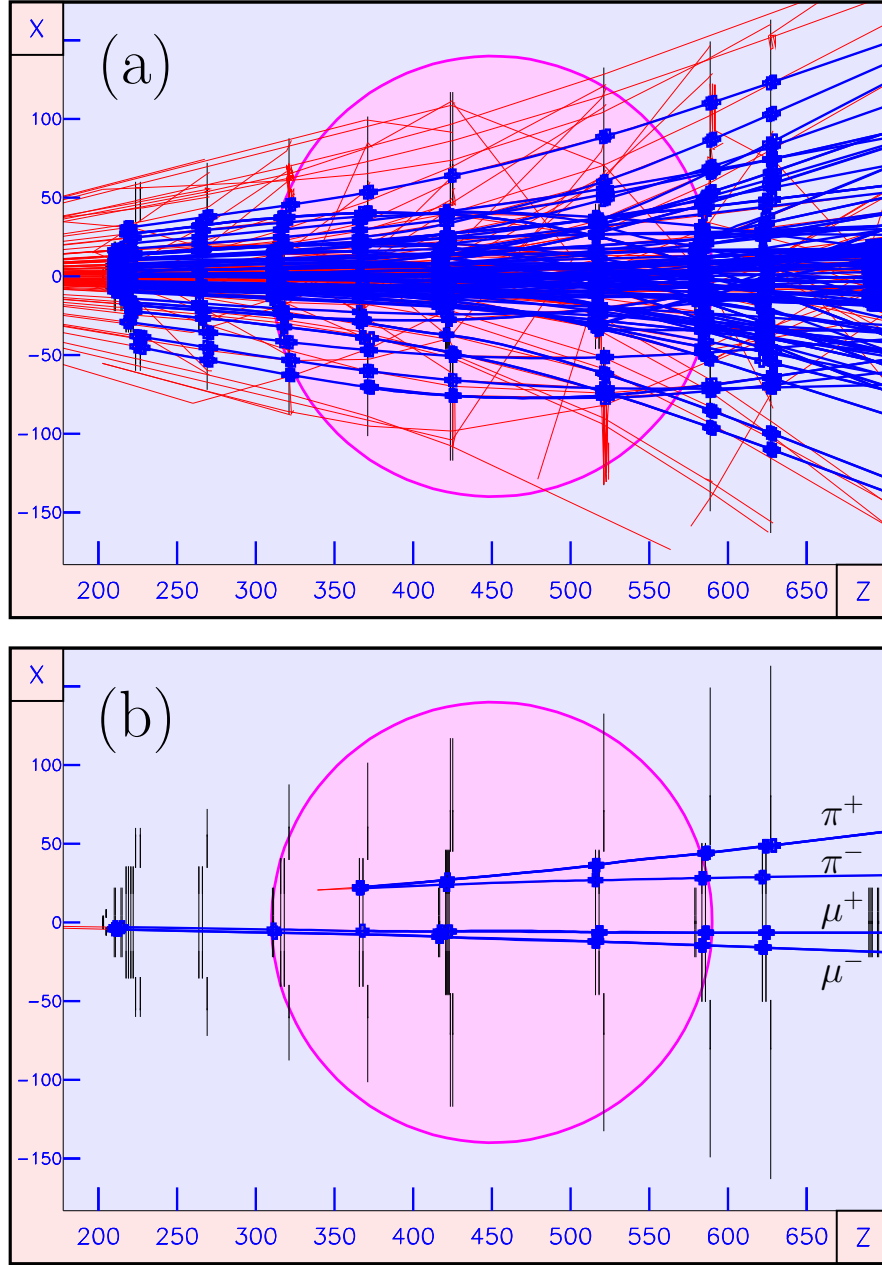


Figure 40: (a) Display of a simulated event with one interaction containing the golden B decay and six superimposed inelastic interactions, focussed on the magnet area, where the pole shoe of the magnet is indicated by the large circle [65]. Both the Monte Carlo tracks (light grey) and the reconstructed tracks (thick dark lines) are shown (reconstructed hit points denoted by crosses). (b) Same event, with the display restricted to particles from the golden B decay.

5 Fitting of Particle Trajectories

After pattern recognition has done its work, the detector hits are separated into sets each of which, ideally, contains manifestations of one specific particle. It is then the task of the track fit to evaluate the track parameters and thus the kinematical properties of the particle with optimal precision. Even if the pattern recognition itself is already providing track parameters and covariance matrices to some degree, obtained for example by means of the Kalman filter, it will in general be left to a final track fit to take all necessary effects into account which are often neglected at the track finding stage because they are costly to apply under the full combinatorics of pattern recognition.

5.1 Random Perturbations

In the easiest case, track parameters could be derived from the measurements by applying the least squares fit formulas from eq. 4 and 5 in sec. 2.4.1. In realistic applications, the problem is usually more involved because of the way the trajectory of the particle is influenced by random perturbations that dilute the information content of the measurements, most commonly multiple scattering and ionization or radiative energy loss. Their influence is schematically displayed in fig. 41. One can interpret the diagram in such a way that, from step to step, the measurements, labelled on the right side, improve the degree of amount of information about the kinematical properties of the particle, while the perturbations labelled on the left side reduce it.

5.2 Treatment of Multiple Scattering

Multiple scattering occurs through the elastic scattering of charged particles in the Coulomb field of the nuclei in the detector material. Since the nuclei are usually much heavier than the traversing particles, the absolute momentum of the latter remains unaffected, while the direction is changed. If the longitudinal extension of the traversed material block can be neglected (this is normally referred to as *thin scatterer approximation*), only track parameters related to particle direction are affected directly, for example the track slopes $t_x = \tan \theta_x$ and $t_y = \tan \theta_y$ introduced in section 2.3.1. The stochastic nature of multiple scattering is that of a Markov process.

The distribution of the deflection angle follows a bell-like shape, though it cannot be accurately described by a Gaussian because of its pronounced tails. The variance of the projected multiple scattering angle is calculated within Molière theory [70, 71, 72] as

$$C_{MS} = \left(\frac{13.6 \text{ MeV}}{\beta pc} \right)^2 t [1 + 0.038 \ln t]^2 \quad (45)$$

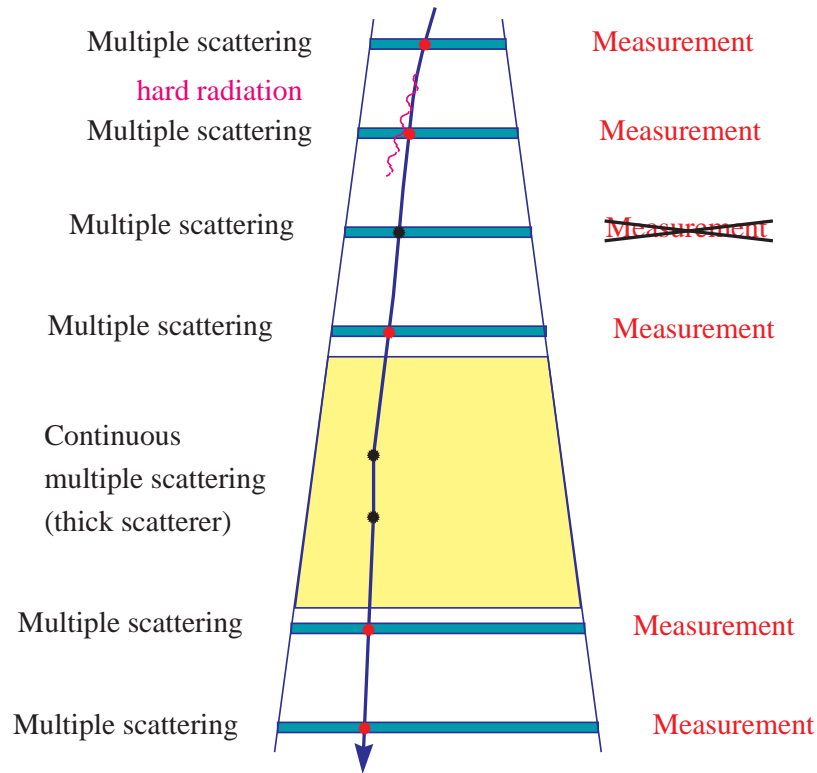


Figure 41: Schematic view of the information flow in the track fit. Elements shown are measurements (hits) in the tracking layers, in one case a missing hit which is e.g. not found by the pattern recognition procedure, and random perturbations like multiple scattering and photon radiation.

where t is the traversed path length in terms of radiation lengths x_R , usually called *radiation thickness*. (While the radiation length is frequently abbreviated as x_0 in the literature, the symbol x_R is used here instead to avoid confusion with other uses of x_0 throughout this article.) For a planar object arranged in a plane vertical to the z axis, the radiation thickness along z is given by

$$\tilde{t} = \int \frac{dz}{x_R(z)} \quad (46)$$

Taking the track inclination against the z axis into account, one obtains the effective radiation thickness

$$t = \tilde{t} \sqrt{1 + t_x^2 + t_y^2} \quad (47)$$

so that the final formula becomes (assuming $\beta \approx 1$)

$$C_{MS} = \left(\frac{13.6 \text{ MeV}}{pc} \right)^2 \sqrt{1 + t_x^2 + t_y^2} \tilde{t} \left[1 + 0.038 \ln \sqrt{1 + t_x^2 + t_y^2} \tilde{t} \right]^2 \quad (48)$$

In general, multiple scattering could be treated in the track fit by expressing the angular uncertainty of each thin scatterer as an additional contribution to the error of each affected measurement. Since a multiple scattering deflection will influence all downstream measurement errors in a correlated way, this introduces artificial correlations into the hitherto uncorrelated measurements, so that the matrix \mathbf{V} in section 2.4.1 is no longer diagonal. Evaluation of eq. 4 requires then inversion of non-trivial matrices whose dimension is not only the number of parameters but the number of measurements. Straight-forward solutions of this problem have been devised [73], which intrinsically treat all multiple scattering angles as free parameters. In many practical situations however, where the number of parameters may be five and the number of measurements perhaps as large as 70, this can lead to serious problems.

The generally accepted solution for the above problem is provided by the Kalman filter technique. The multiple scattering dilution is added as *process noise* (represented by the matrix Q_k in the transport equation, eq. 8) at the very position in the trajectory where it originates. The Kalman filter normally proceeds in the inverse flight direction along the path of the particle and takes the influences illustrated in fig. 41 into account. Mathematically, the result will be identical to a straight-forward least squares fit as described in the previous paragraph, but the detailed procedure avoids handling of huge matrices.

In Kalman filter language, the resulting covariance matrix contribution for thin scatterers is

$$\begin{aligned} \text{cov}(t_x, t_x) &= (1 + t_x^2)(1 + t_x^2 + t_y^2)C_{MS} \\ \text{cov}(t_y, t_y) &= (1 + t_y^2)(1 + t_x^2 + t_y^2)C_{MS} \\ \text{cov}(t_x, t_y) &= t_x t_y (1 + t_x^2 + t_y^2)C_{MS} \end{aligned} \quad (49)$$

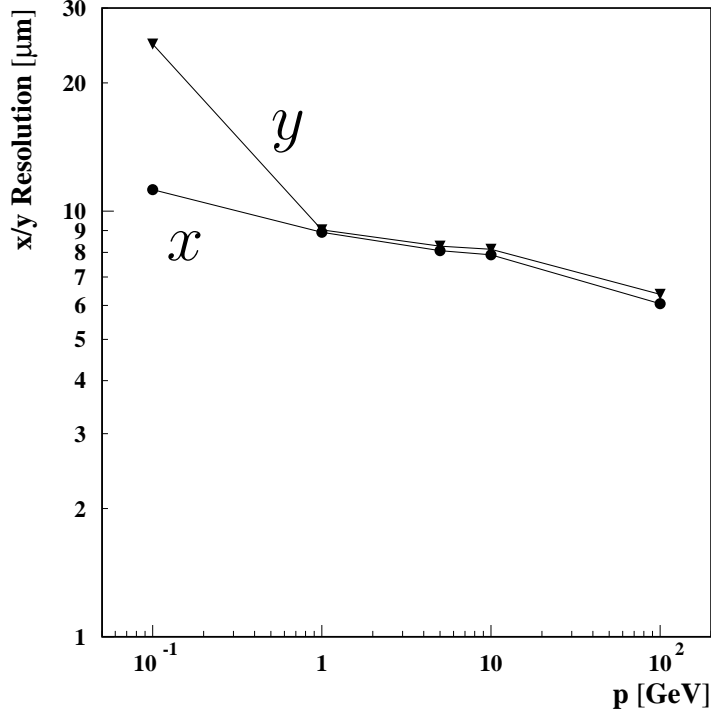


Figure 42: Impact parameter resolution at the first track point, separately for the coordinate in the bending plane (x , circles) and the non-bending plane (y , triangles).

(These and related formulas and their derivation can be found in [74]).

It may be interesting to see how such a fit works in practice. In the following, results of a study are shown which has been performed on basis of simulated events in the HERA-B geometry (fig. 8), applying a Kalman filter based track fit to the simulated hits [75]. This kind of geometry is typical for modern forward spectrometers, and generally similar to COMPASS [76] or the planned LHCb [77] and bTEV [78]. The study was based on detector design resolutions and not intended to make quantitative statements on the actual spectrometer performance, but to provide insight into the effects of combining various different detector types, the sizable number of hits per track, and the considerable amounts of material in the tracking area that make an accurate treatment of multiple scattering essential.

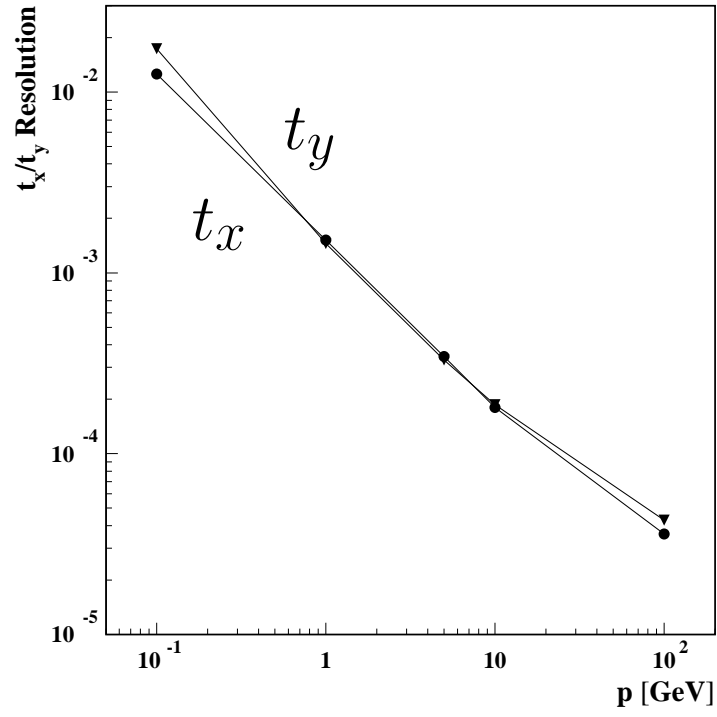


Figure 43: Resolution of the slope parameters $t_x = \tan \theta_x$ (circles) and $t_y = \tan \theta_y$ (triangles).

5.2.1 Impact parameter and angular resolutions

The visible track parameter resolution was obtained by calculating the *track parameter residual* for each track using the Monte Carlo truth, and applying a Gaussian fit to the distribution. (The term *visible* is used to distinguish this resolution from the one estimated by the fit.) The impact parameter resolution for tracks passing the Silicon micro-vertex detector and the outer tracker as a function of momentum is shown in fig. 42. Since this impact parameter is defined with respect to the position of the first hit of the track counting from the interaction point, the resolution is governed by the error of the first coordinate and only weakly dependent on momentum. Multiple scattering acts like a filter which dilutes the information from the following layers, only at higher momentum their contribution to the resolution at the first point becomes visible.

Since the vertex detector measurement accuracy is approximately isotropic, horizontal and vertical resolution are almost identical, the deviation at $p = 100$ MeV is explained by the fact that the strips in the first vertex detector layer are oriented almost parallel to the y axis. The resolution of track slopes is shown in fig. 43 and turns out to be dominated by the pronounced $\propto 1/p$ behaviour expected in a multiple scattering-dominated regime. At high momentum, the onset of coordinate resolution effects appears to be just visible, where the slightly better resolution of the horizontal slope (t_x) may be due to the dominantly vertical orientation (parallel to y) of the wires in the main tracking system.

The impact parameter resolution given above should not be confused with the quantities relevant for physics performance where assignment to vertices is important. In the latter case, the track parameters must be extrapolated from the first track point to the interaction area. With extrapolation distances of typically $\mathcal{O}(10\text{ cm})$, the resolution of the *extrapolated* impact parameters will generally be fully dominated by the *angular* resolution rather than the impact parameter resolution at the first point.

5.2.2 Momentum resolution

A very central design issue in spectrometers is resolution of momentum, since it determines the rejection power against background in particle spectrometry. The relative momentum resolution, labelled dp/p , as a function of momentum is shown in fig. 44 for particles traversing the areas *SI*, *MC* and *PC* of the spectrometer (see fig. 8 for definition) in the polar angle area $0.1 < \theta < 0.15$. The circle symbols show the relative momentum resolution that results with multiple scattering switched off in the simulation, leading to a strictly linear dependence on p . This behaviour is expected since the resolution is then only determined by the coordinate resolution and the geometrical layout of the spectrometer - size and number of layers - that provides the leverage for momentum measurement together with the magnetic field. The result reflects the fact that the curvature κ ,

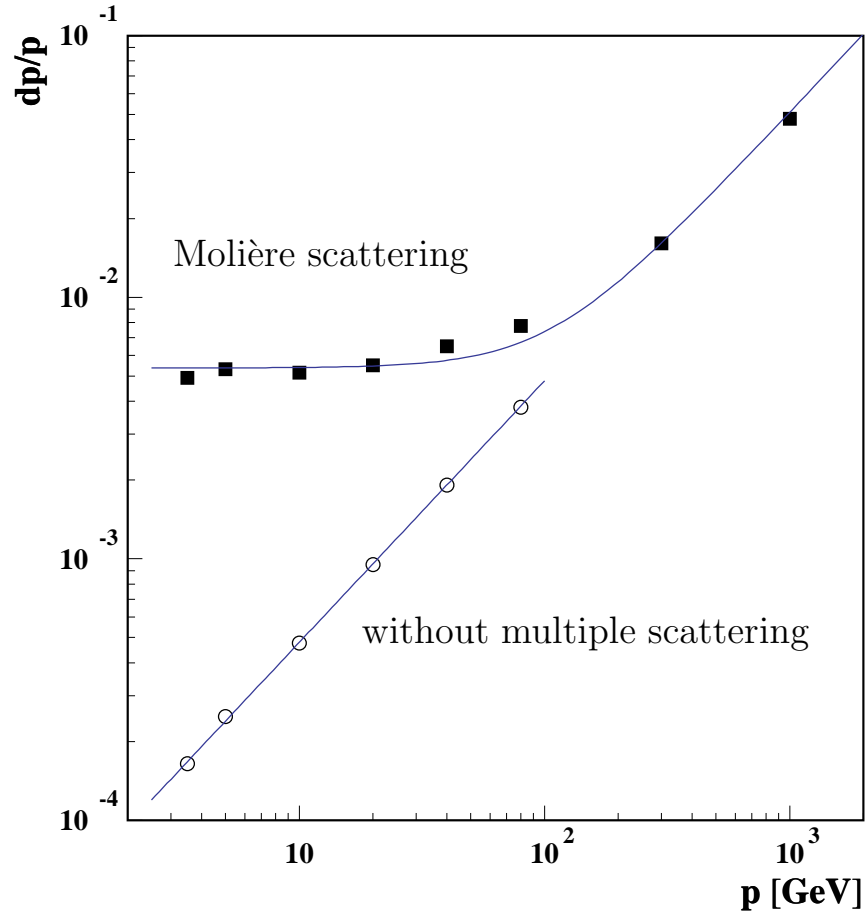


Figure 44: Visible momentum resolution (filled squares) for simulated muons (filled squares) together with the fit of the parametrization described in the text (upper solid line). In the lower part, the open squares show the visible resolution with multiple scattering switched off in the simulation, with a similar function fitted. The open circles show the pure coordinate resolution as estimated by the track fit, with a linear function fitted to it.

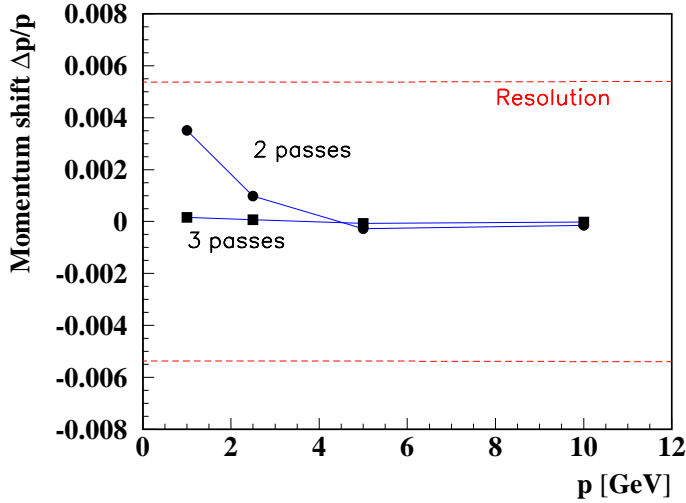


Figure 45: Residual of the momentum parameter (see eq. 23) normalized to the momentum itself for two and three passes of the fit. The relative momentum resolution is indicated by the dashed lines for comparison.

which is the inverse of the radius of curvature, can be measured with a precision that is independent of its actual value, hence $\delta\kappa=\text{const}$. On the other hand the curvature is inversely proportional to the momentum, so that $dp/p \propto p$. In presence of multiple scattering, the resolution shows a multiple scattering-dominated regime below momenta of ≈ 50 GeV, and a transition into a linear rise at high momentum. Superimposed is a fit with a constant and a linear resolution term added in quadrature. This parametrization, which corresponds to a commonly used function introduced by Gluckstern [7] for an even spacing of tracking stations does not fit the visible resolution very well in the momentum mid-range, which can be attributed to the uneven distribution of measurements, resolutions, material and magnetic field strength in the spectrometer.

5.2.3 Effects of fit non-linearity

The presence of the inhomogeneous magnetic field introduces particular effects of non-linearity into the fitting problem. The least squares fit technique, which the Kalman filter is built on, can still be applied, with the transport matrices now obtained as derivatives of the transport function. As already noted in sec. 2.4.1, the optimal properties of the least squares technique are still retained on the condition that the derivatives are taken at the position of the final trajectory. Since this is initially not necessarily the case, the fit must be repeated iteratively until the procedure converges.

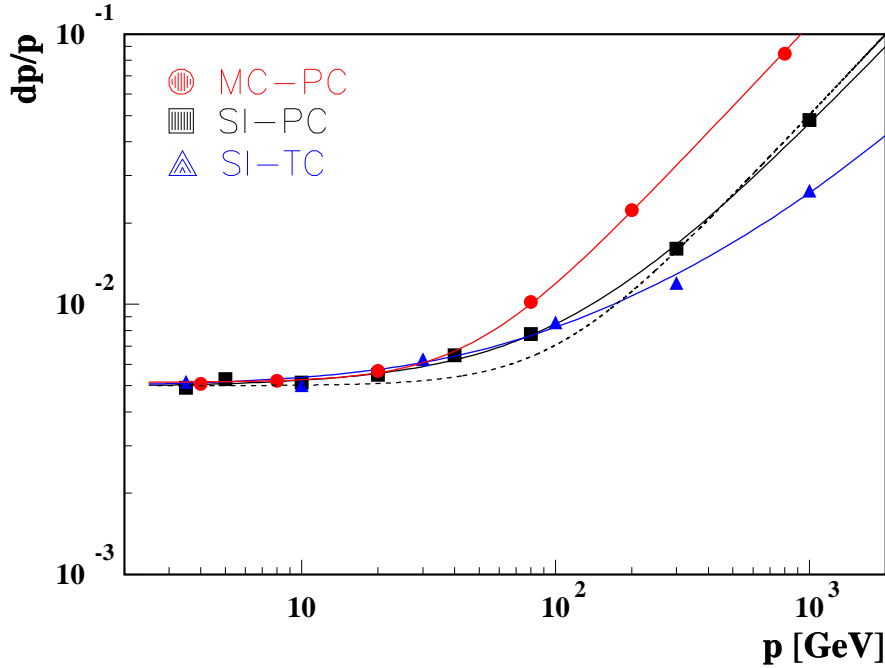


Figure 46: Relative momentum resolution in the MC-PC (circles), SI-PC (squares) and SI-TC (triangles) spectrometer ranges, together with the fits described in the text. The dashed line is the upper fit in fig. 44.

The practical implications of non-linearity are visible in fig. 45, which shows the mean relative deviation of the reconstructed from the true momentum value. Small systematic shifts of reconstructed momentum are observed for momentum below 5 GeV with two fit passes applied. These shifts reflect the convergence behaviour of the fit due to non-linearity. They are found to be virtually removed when a third pass is applied.

5.2.4 Contributions of different parts of the spectrometer

For understanding detector design, it is also important to investigate how much different parts of the spectrometer contribute to the momentum measurement. In the HERA-B geometry (fig.8), the tracking system is grouped into the vertex detector (SI), the chambers within the magnet (MC), the chambers just behind the magnet (PC) and the so-called *trigger chambers* (TC), which are separated from the PC part by the ring-imaging Čerenkov detector (RICH). In order to separate the contributions of the different spectrometer parts, the range of the fit was modified by omitting the vertex detector hits (labelled MC-PC range) and by adding the hits from the tracking chambers at the end of the main tracking

system (SI–TC range). The resulting momentum resolutions are displayed in fig. 46. It turns out that without including the vertex detector (MC–PC), the momentum resolution is well described by a constant and a linear term added in quadrature. In the regime of linear rise, the poorer coordinate resolution is reflected in comparison to the system including the vertex detector. When the fit on the other hand is extended into the “TC” region which is mainly designed to support the trigger (SI–TC), these additional measurements with their huge lever arm are expected to improve the coordinate contribution of the resolution. Such an improvement is visible in fig. 46 for $p \geq 100 \text{ GeV}$, where it is hardly relevant for the physics scope of the experiment. A third term proportional to the square-root of the momentum had to be added in quadrature to fit the resolution for the latter two ranges.

5.2.5 Parameter covariance matrix estimation

A very important task of the track fit is the quantification of the covariance matrix of the estimated track parameters. The reliability of parameter error estimation can be studied by investigating distributions of *normalized parameter residuals* (see eq. 23 in sec. 2.5.5), which use the estimated error for normalization. In the example at hand, the resulting pull distributions are shown in fig. 47, where unbiased fits with a Gaussian function are superimposed. Distortions of the parameter estimates would show up as deviations of the mean values from zero, which are however not present in this case. The Gaussian cores of the pulls agree in all cases with unity width, indicating a reliable estimate of the covariance matrix. One should note that only mean value and variance of the pull distribution are indicators of the quality of the estimate. The actual *shape* of the distribution, e.g. whether it is Gaussian or not, reflects the underlying structure of the problem, as will be more clearly visible in the next section.

5.2.6 Goodness of fit

Since the Kalman filter is mathematically equivalent to a least-squares estimator, the sum of the *filtered* χ^2 contributions will follow a χ^2 distribution, provided that the random variables entering into the fit have Gaussian distributions. In this case the χ^2 *probability*

$$P_{\chi^2} = \int_{-\infty}^{\chi^2} f(\tilde{\chi}^2) d\tilde{\chi}^2$$

where $f(\tilde{\chi}^2)$ is the standard χ^2 distribution for the appropriate number of degrees of freedom, should be evenly distributed between 0 and 1. (P_{χ^2} is often called *confidence level*.) This prerequisite is not strictly fulfilled in case of Molière scattering, so that deviations are to be expected. These effects have potentially large influence in modern radiation hard drift chambers, where the drift cells

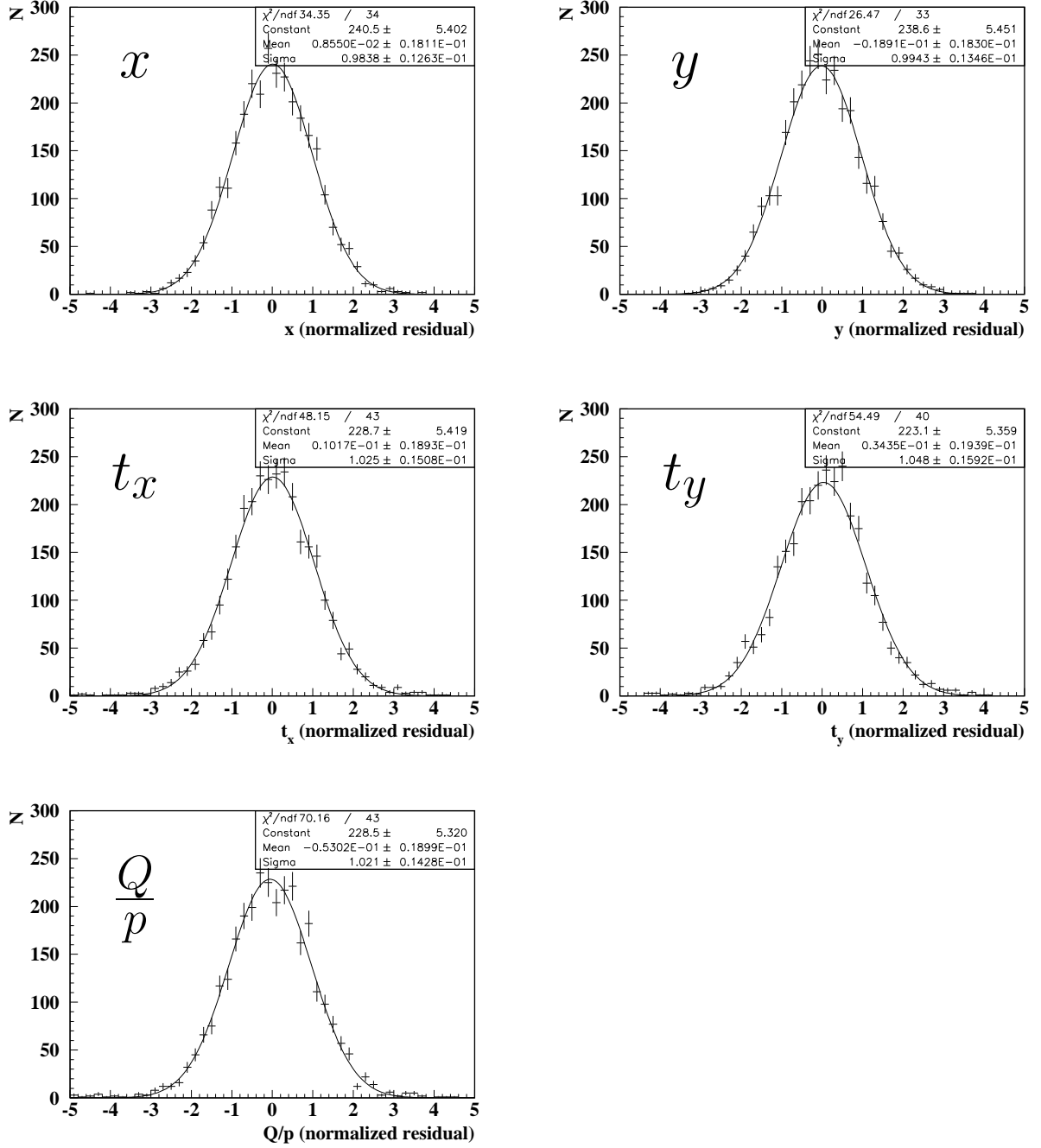


Figure 47: Normalized parameter residual distributions for muons of 10 GeV, based on 3000 simulated tracks.

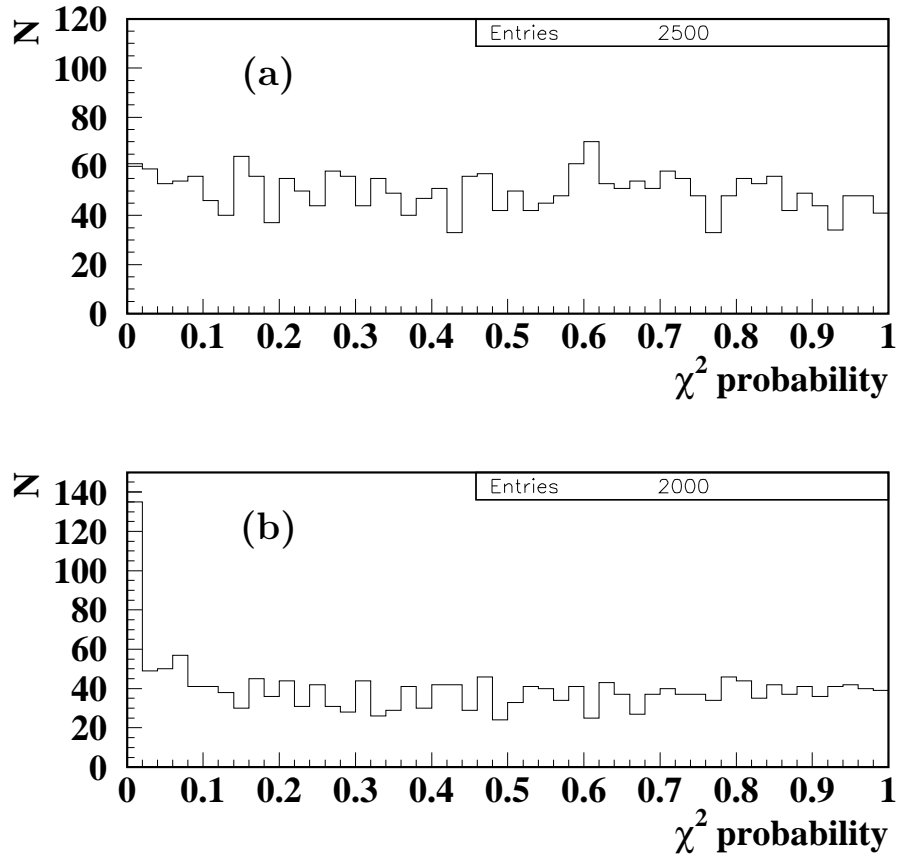


Figure 48: Distributions of χ^2 probability (confidence level) for the track fit for a 10 GeV particle, (a) with Gaussian form of multiple scattering, (b) with Molière scattering.

are enclosed in a multitude of small gas volumes and a considerable amount of material is introduced into the tracking area.

Figure 48 compares the distribution of the χ^2 probability for the Gaussian form of multiple scattering (a) and Molière scattering (b). The peak at small probabilities in (b) obviously does not indicate a bad behaviour of the fit, but instead shows the inadequateness of the χ^2 test with non-Gaussian random variables. The probability distribution for various momentum values is displayed in fig. 49. The increasing prominence of the peak at low probability is clearly seen with decreasing momentum. Small χ^2 probability does not necessarily imply a bad estimation of the parameters, hence special care is required when a χ^2 cut is to be used to eliminate improperly reconstructed tracks.

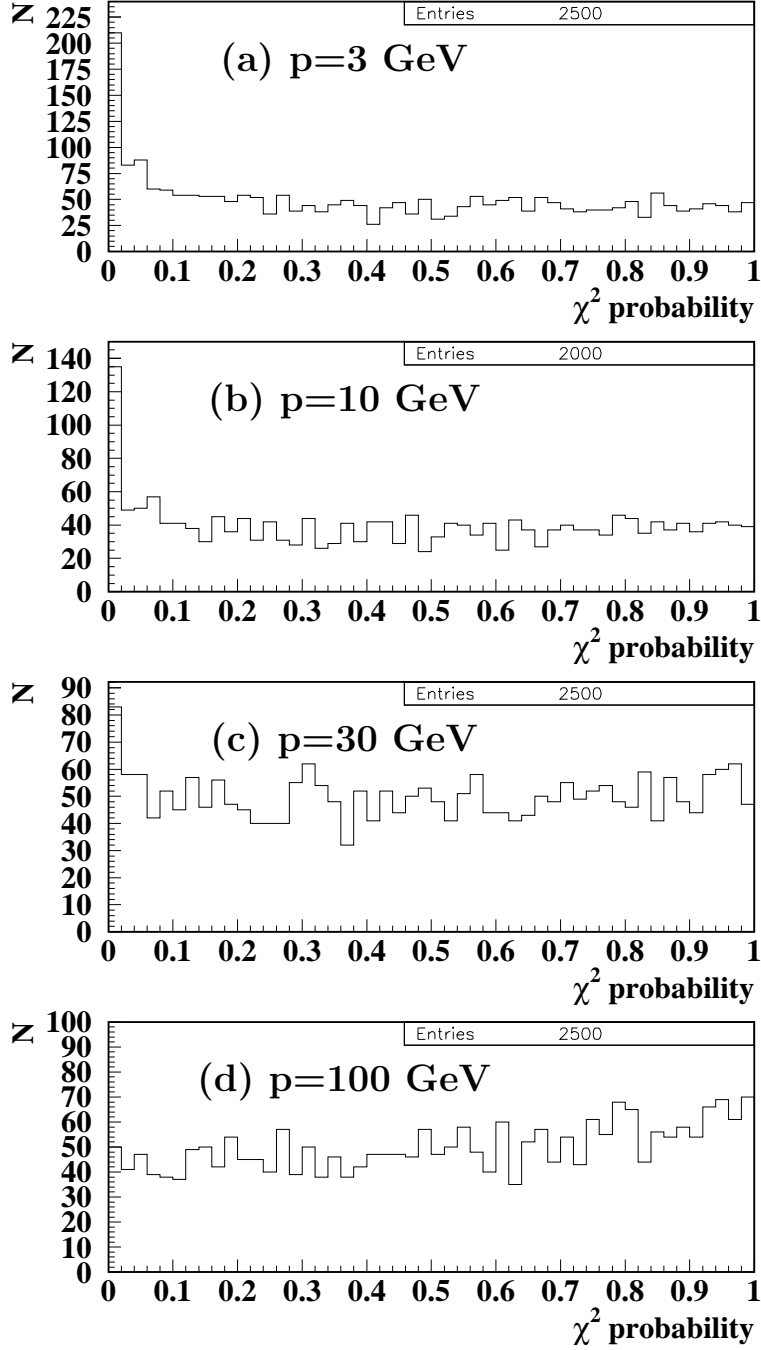


Figure 49: Distribution of χ^2 probabilities as a function of momentum, (a) 3.5 GeV, (b) 10 GeV, (c) 30 GeV and (d) 100 GeV.

5.3 Treatment of Ionization Energy Loss And Radiation

5.3.1 Ionisation energy loss

For minimal ionizing particles in the GeV energy range, energy loss due to ionization within the tracking system depends in good approximation only on the amount of material that is traversed. In this case, it is not the radiation thickness (as defined in eq. 46), but the geometrical thickness multiplied by the mass density of the material that is relevant. Since the energy loss depends only weakly on the energy itself in this range, the effect will become most noticeable for low momentum particles. This behaviour is illustrated in fig. 50, which shows the normalized residual of the momentum parameter Q/p for μ^+ particles of 3.5 and 10 GeV with ionization energy loss simulation turned on. The residual distributions are shifted towards positive values of Q/p , reflecting an underestimation of the energy, which is caused by the ionization energy loss, in particular upstream of the magnet. The visible shift corresponds to an energy loss of 12 MeV. On the other hand, the width of the residual distributions is not significantly increased, which in the 10 GeV case can directly be seen by comparing with fig. 47.

A correction can be applied in each filter step if the dE/dx of the particle in the material is known, since

$$E_{after} = E_{before} - (dE/dx)_{ion} \cdot \ell \quad (50)$$

where ℓ is the traversed thickness of the material. This requires in general the knowledge of the particle mass. Since ionization energy loss will be most notable for small particle energies where the resolution is governed by multiple scattering, no correction to the momentum error has been applied. The bottom part of fig. 50 displays the same normalized residuals with the energy loss correction applied. The bias of the momentum estimate is successfully eliminated by the correction.

5.3.2 Radiative energy loss

The corrections discussed up to now are usually sufficient for minimum ionizing particles. For electrons⁶ however, the situation is more complicated since above the *critical energy*, which is of the order of MeV, these particles lose more energy through radiation of photons than through ionization when they traverse material. This process is also of a more notably stochastic nature than ionization energy loss, as considerable fractions of the electron energy can be transferred to the photon. Modern radiation-hard detectors as e.g. those under construction for the LHC are confronted with this problem to a much higher degree than traditional detectors, because of the significant amount of material in the tracking system, which can easily exceed 50% of a radiation length.

⁶in this section the term *electron* should be interpreted to imply *positron* as well

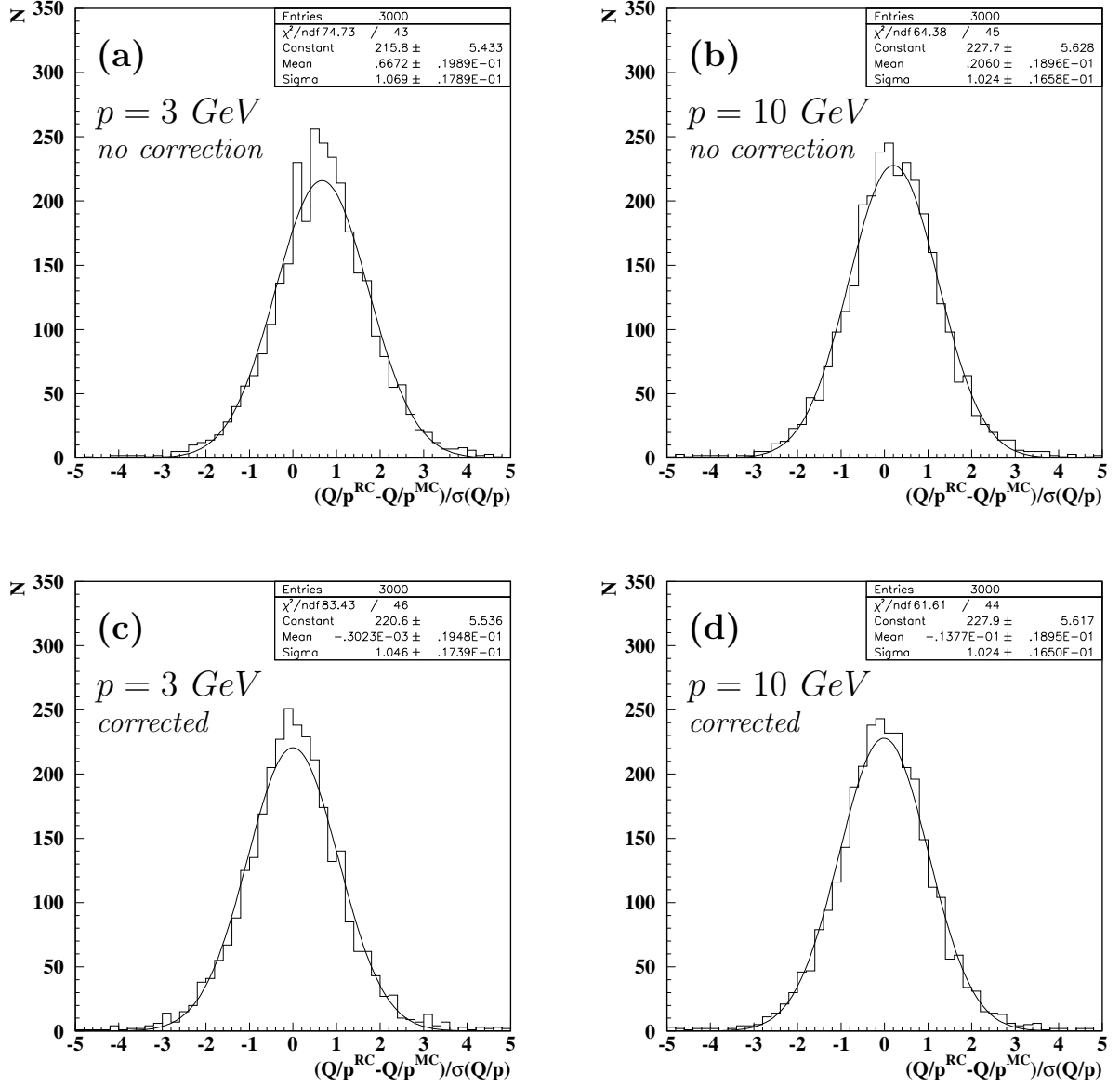


Figure 50: Pull distribution of the momentum parameter for μ^+ particles of 3 GeV (a,c) and 10 GeV (b,d). The upper pictures show the effect of dE/dx if no correction is applied. The lower plots show the same when the correction is applied in the fit.

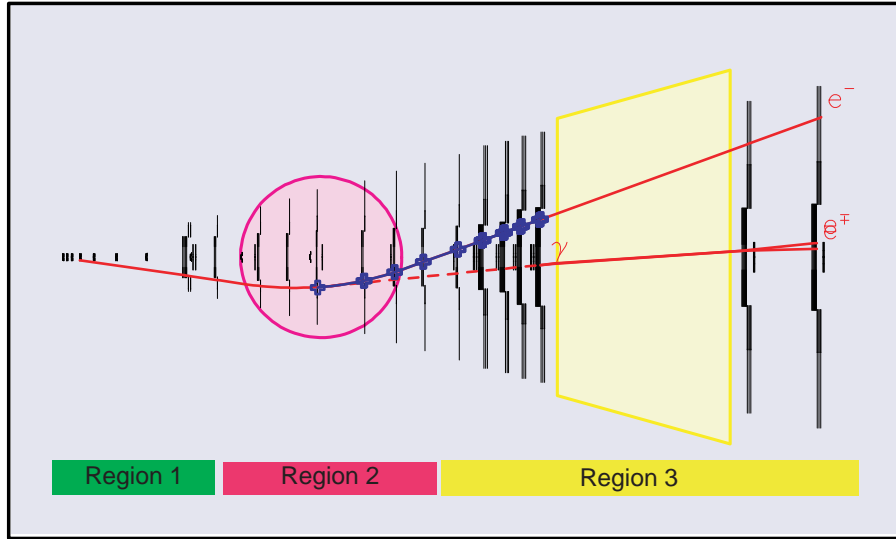


Figure 51: Regions 1–3 for classifying radiative energy loss illustrated in the geometry of the HERA-B spectrometer. The simulated geometry differs in some details from the one in fig. 8. Also the trajectory of a simulated electron is shown, which radiates a photon within the magnet that converts into a e^+e^- pair further downstream.

For the relevance of photon radiation on measurement of the electron, three cases have to be distinguished regarding the range where the radiation occurs (indicated as regions 1–3 in fig. 51):

Region 1: between interaction point and spectrometer magnet If the point of origin of the particle is not yet within the magnetic field – as is typical for fixed-target setups rather than for collider detectors – radiation will not change the electron trajectory and thus not interfere with the quality of the fit; however, the spectrometer will only measure the remaining momentum of the electron after the radiation.

Region 2: within the magnetic field In this case, the curvature of the trajectory changes because of the radiation, which means that the energy change is – in principle – measurable. Ignoring the radiation in the fit will lead to a bad description of the trajectory and to distortions of the parameter estimates.

Region 3: beyond the magnetic field If the electron loses energy downstream of the magnet, this will have no influence on the momentum measurement in the spectrometer. However, pair creation from radiated photons may lead to accompanying particles that can disturb pattern recognition in the downstream area.

The dilution due to energy loss of electrons and positrons through emission of electromagnetic radiation can be treated by the method by Stampfer et al. [79]. According to the Bethe-Heitler equation [80], this energy loss is described by

$$\left(\frac{dE}{dx}\right)_{rad} = \frac{E}{x_R} \quad (51)$$

where x_R is the radiation length of the traversed material (see section 5.2). This leads to the relation

$$\left\langle \frac{E_{after}}{E_{before}} \right\rangle = e^{-t} \quad (52)$$

where t is the traversed distance measured in radiation lengths as defined before. For a track propagation which follows the track opposite to its physical movement, one obtains on average

$$\left(\frac{Q}{p}\right)' = \frac{Q}{p} + \Delta \left(\frac{Q}{p}\right) = \frac{Q}{p} - \frac{Q}{p} \frac{E_{before} - E_{after}}{E_{before}} = \frac{Q}{p} e^{-t} \quad (53)$$

The contribution to the propagated covariance matrix emerges as

$$\Delta_{cov} \left(\frac{Q}{p}, \frac{Q}{p} \right) = \left(\frac{Q}{p} \right)^2 \left(e^{-t \frac{\ln 3}{\ln 2}} - e^{-2t} \right) \quad (54)$$

This contribution can be included into the Kalman filter process noise as introduced in eq. 8.

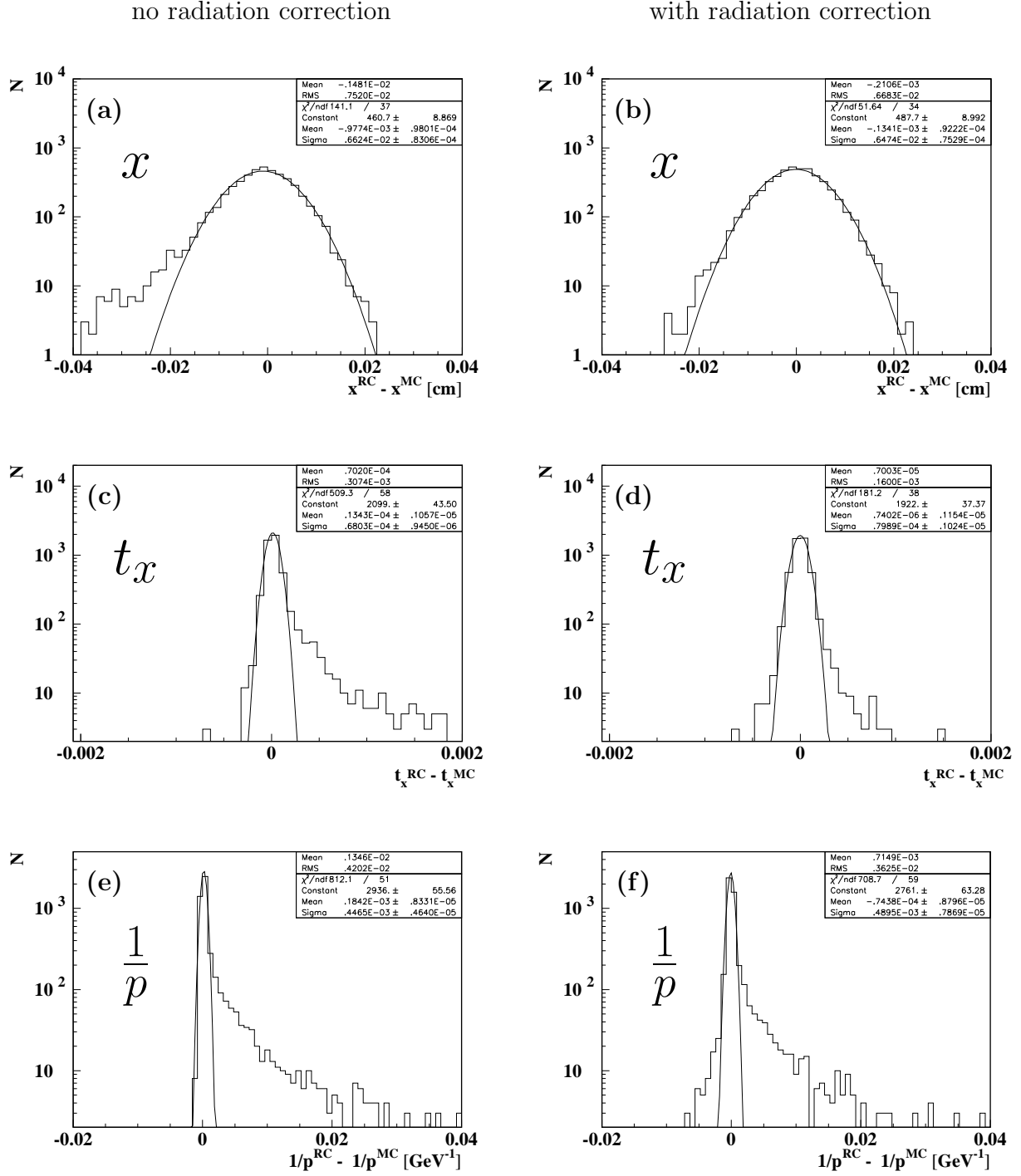


Figure 52: Distributions of parameter residuals for electrons of 100 GeV based on 5000 tracks, where x is the impact parameter in the bending plane, $t_x = \tan \theta_x$ is the corresponding track slope, and $1/p$ the inverse momentum. The track fit was applied to all hits within the spectrometer magnet (region 2 in fig. 51).

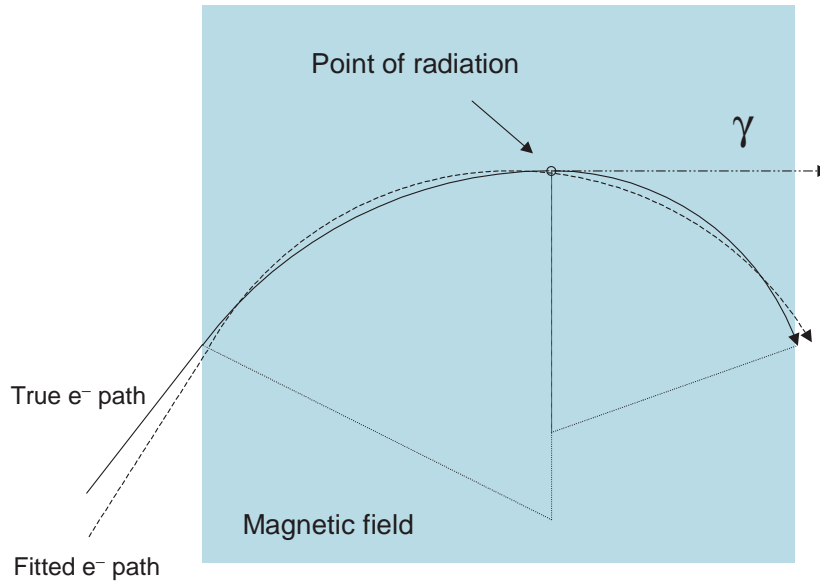


Figure 53: Illustration of how radiation within the magnetic field can affect the estimate of the fitted track slope. The magnetic field vector is pointing into the drawing plane. The electron, whose true path is shown by a solid line, emits a photon, which leads to an increase of curvature for the subsequent part of its trajectory within the magnet. This is illustrated by the curvature radii of the helices as dotted lines. The fitted trajectory (dashed line) assumes a single curvature, which leads to an overestimation of the initial slope of the track. The curvatures drawn are intentionally exaggerated.

| Radiation correction mode | Fraction of fits within momentum deviation | |
|---------------------------|--|----------------------------|
| | $-0.1 < \delta p/p < +0.1$ | $-0.2 < \delta p/p < +0.2$ |
| none | 0.566 ± 0.004 | 0.678 ± 0.003 |
| within magnet | 0.635 ± 0.003 | 0.728 ± 0.003 |
| within full spectrometer | 0.321 ± 0.003 | 0.786 ± 0.002 |

Table 3: Fraction of fits within given limits of momentum deviation, for three variants of radiation correction

5.3.3 Radiation energy loss correction within the magnetic field

Energy loss through radiation can not only interfere with the momentum measurement, but may also affect other track parameters. This is shown in figs. 52a,c,e which display the residuals of the parameters x , t_x and $1/p$ for electrons produced with 100 GeV momentum, where the fit was restricted to the magnet area (MC). Without bremsstrahlung correction, the track slope estimate t_x shows a tail towards overestimated values, which is reflected in an underestimation of the corresponding impact parameter, x . The explanation for this effect is illustrated in fig. 53 which for simplicity assumes a homogeneous field: the curvature of the electron track is abruptly increased beyond the point of radiation. Fitting the track with a constant momentum leads to an intermediate curvature resulting in a shift in the measured initial track slope.

The residual distribution of the momentum parameter, $1/p$, displays a tail towards higher values, corresponding to a mean momentum shift of $\approx 13\%$.

Also the parameter errors are underestimated, which is evident from the normalized residuals in figs. 54a,c,e (uncorrected case), where the widths of the t_x and Q/p pull distributions are significantly enlarged.

Figures 52b,d,f show the result with the radiation correction applied in the fit. One can see that the tails in the parameter estimates of x and t_x are far less pronounced, and the bias in the impact parameter and track slope is considerably reduced. Also the distortion of the mean reconstructed inverse momentum $\delta(1/p) \approx \delta p/p^2$ is reduced from $1.3 \cdot 10^{-3} \text{ GeV}^{-1}$ to $7 \cdot 10^{-4} \text{ GeV}^{-1}$, and the standard deviation (RMS width) of the parameter estimates is reduced by 11% (x), 48% (t_x) and 14% (Q/p), respectively. Moreover, the radiation correction brings the RMS widths of the pull distributions close to unity (figs. 54b,d,f), which indicates a reliable covariance matrix estimate. The fit probability distribution is shown in fig. 55. It reflects a non- χ^2 type distribution of the goodness-of-fit, which is expected since the radiation of bremsstrahlung introduces a strongly non-Gaussian random perturbation.

The situation is different if one attempts to extend the radiation correction to the full tracking system including regions 1 and 3 which are outside of the magnetic field, most notably the vertex detector whose material causes a significant energy loss for electrons. Outside of the magnetic field, however, the trajectory

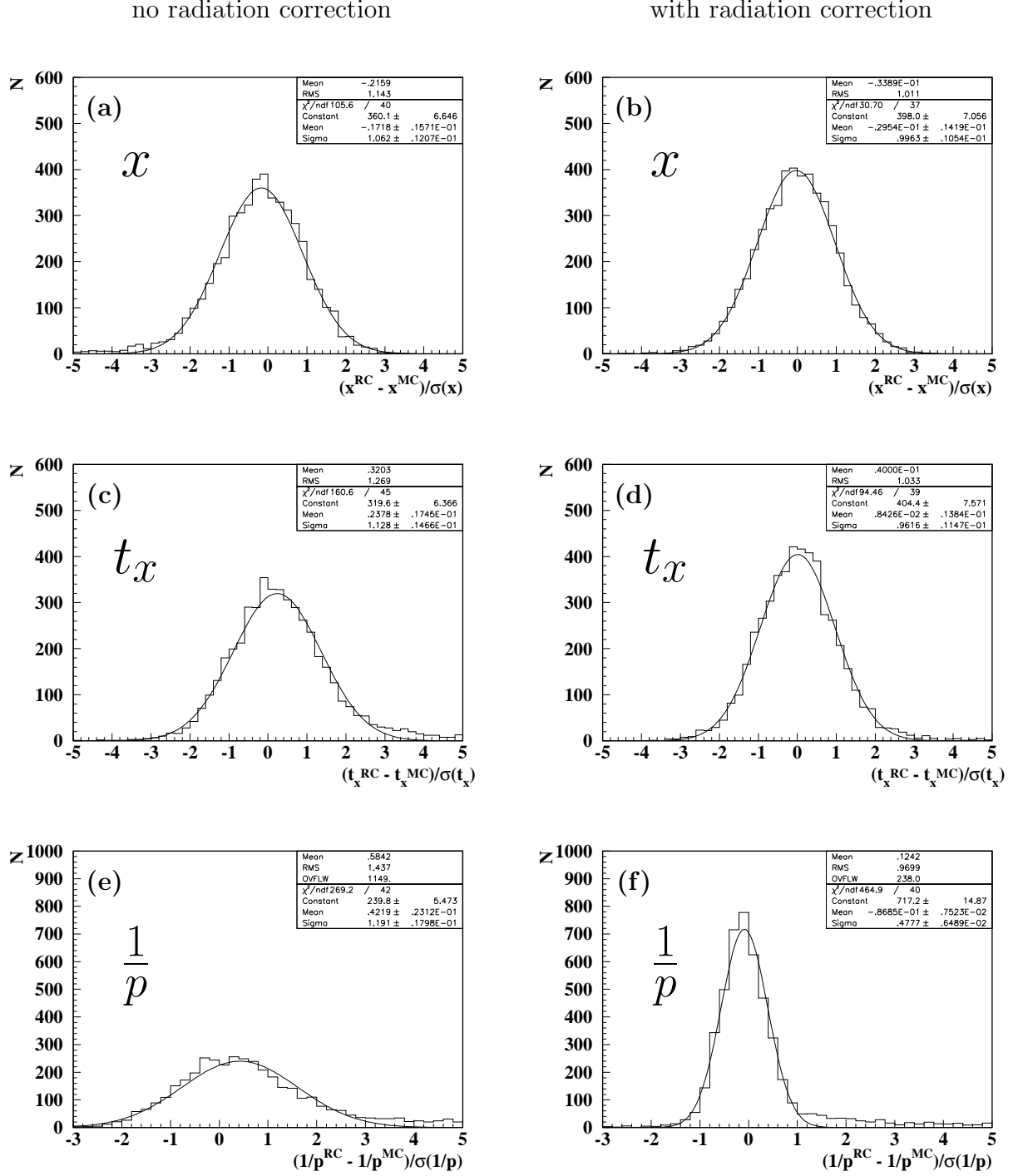


Figure 54: Distribution of normalized parameter residuals (pulls) for electrons of 100 GeV based on 5000 tracks, where x is the impact parameter in the bending plane, $t_x = \tan \theta_x$ is the corresponding track slope, and $1/p$ the inverse momentum. The track fit was applied to all hits within the spectrometer magnet (region 2 in fig. 51).

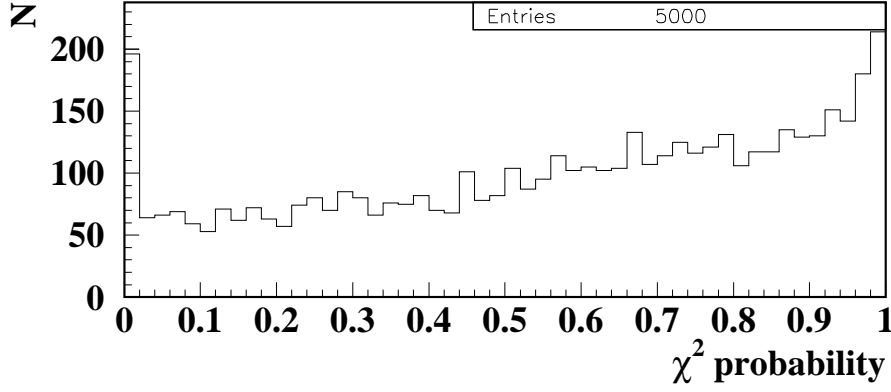


Figure 55: Distribution of χ^2 probability from the track fit of 100 GeV electrons in the main tracker with radiation correction. The non-Gaussian distribution of the radiated energy leads a non-flat probability distribution with a sharp peak near zero.

full radiation correction

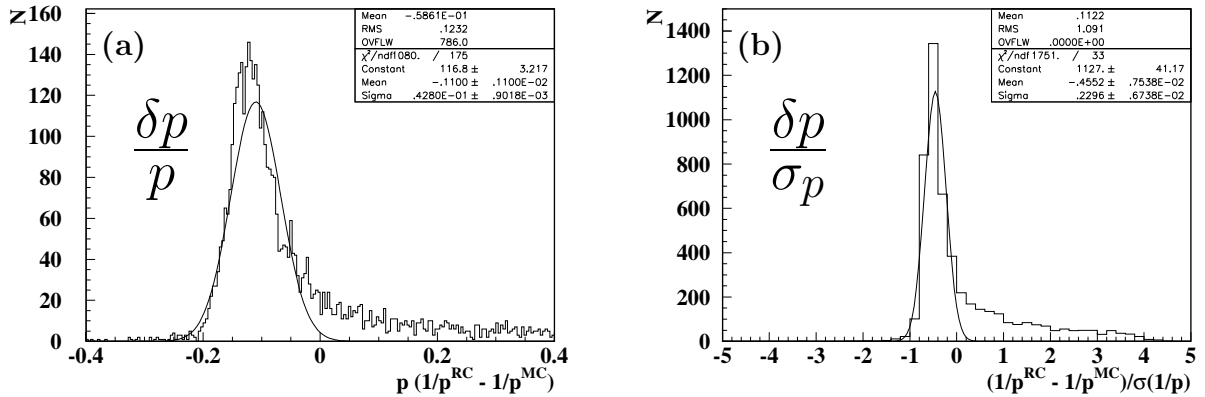


Figure 56: Distribution of $1/p$ parameter residuals multiplied by p as measure for relative momentum deviation (a), and of normalized $1/p$ residuals (b) for 100 GeV electrons in the full tracker. The radiation correction was applied in the whole tracking system, leading to the shift of the peaks described in the text.

shape is not modified by radiation, which means that the fit will only apply the on-average correction according to the traversed radiation thickness. This can lead to bizarre results as seen in fig. 56, which shows the distribution of the $1/p$ parameter residual multiplied by the momentum itself as well as the corresponding pull distribution. The peak has moved away from zero to negative residual values, implying that electrons in the peak obtain an overcorrected energy value. In fig. 56b, the mean value of the pull is near zero, and the RMS width is close to one, indicating that the compensation works correctly in the statistical sense. For intuitive plausibility, however, it is relevant that a large fraction of measurements are in the immediate vicinity of the quoted value. A test of this criterion is shown in table 3, which summarizes the fraction of fits with momentum deviation of within 10% or 20% of the real value for the three correction scenarios. With the 10% criterion, the full spectrometer correction appears worse than even in the uncorrected case, while the correction restricted to the magnet gives the best description in the intuitive sense. In conclusion, the magnet-based correction appears to be provide the best compromise, though this will in general have to be evaluated in each specific application.

5.4 Robust Estimation

The preceding sections have shown how intrinsically non-Gaussian influences, as multiple scattering, or radiative energy loss of electrons, can complicate the estimate of essential kinematic parameters and their interpretation. A fully adequate treatment of profoundly non-Gaussian variables is in general beyond the capabilities of least squares estimation. Likelihood methods, on the other hand, are in principle able to cope with random variables of any distribution, but often cannot be used with as efficient a machinery, in particular when it comes to computation of error matrices.

During the last years, promising concepts have been developed that permit treatment of non-Gaussian random variables, but still allow to use much of the powerful machinery developed with least squares estimation. These methods are called *robust estimation* techniques. One very attractive idea is based on the fact that non-Gaussian distributions can often be approximated as superposition of a limited number of Gaussian distributions [81, 82]. For example, a distribution resembling a Gaussian in the centre, but featuring long tails, as is common with multiple scattering, can be approximated by a sum of a narrow Gaussian distribution and a wide one. If one performs two parallel least squares estimates, each based on one of the Gaussians, the resulting parameter estimates, combined with appropriate weights, will reflect the underlying statistics better than a single estimate with a single Gaussian approximation. Thus, the occurrence of random variables in the tail of the distribution does not pull the estimate as far away as it would with a traditional least square estimator, leading to a more robust behaviour of the fit.

This is the basic idea of the *Gaussian Sum Filter* (GSF) [81, 82, 83, 84, 85], which uses the Kalman filter to incorporate the individual Gaussian components. Upon each occurrence of process noise, the distribution of which is approximated by a sum of N Gaussians, the filter splits into N parallel branches each of which obtains a corresponding weight. In a detector geometry with many scattering elements, this will lead to a repeated multiplication of the number of linear filters to be evaluated. To avoid the explosion of the computing effort, the number of parallel components is limited by *collapsing* or *clustering* components of similar shape. It has been shown that the algorithm can be designed such that the computing effort increases linearly with the maximum number of parallel components (M), and that $M \approx 6 - 8$ already gives good results [84]. In a similar way, radiative energy loss of electrons can be treated by approximating the radiated energy distribution by superposition of several Gaussians [86].

6 Event Reconstruction

After particle tracks have been reconstructed, they form the basis for the reconstruction of the whole event. This will ultimately include particle identification based on dE/dx , time-of-flight, Čerenkov or transition radiation, muon chambers and calorimetry, as well as kinematical reconstruction of composite particles and jets. This article will restrict itself to a brief discussion of vertex reconstruction and kinematical constraints.

6.1 Vertex Pattern Recognition

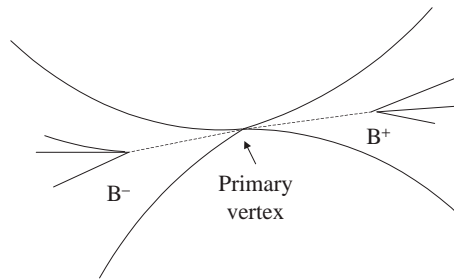


Figure 57: Schematic view of the event structure in an interaction of the type $e^+e^- \rightarrow B^+B^- + X$

The vertex is an essential element of the space-time structure of an interaction. Vertices indicate either the location where an interaction has taken place, for example the primary interaction that is the ultimate origin of all emerging particles, or the place where an unstable particle has decayed. This is illustrated

in fig. 57, which schematically sketches the final state of an interaction with associated production of two beauty mesons, as it can occur for example at a high energy e^+e^- collider. The beauty hadrons, here a B^+ and a B^- , are produced together with accompanying charged particles at the interaction point, travel invisibly for some distance that is, on average, determined by their lifetime and momentum, whereupon they decay into daughter particles. The charged tracks coming from these decays can be used to reconstruct the decay locations of the B mesons as *secondary vertices*⁷. The other tracks, together with the reconstructed B mesons form the *primary vertex*, which indicates the interaction point.

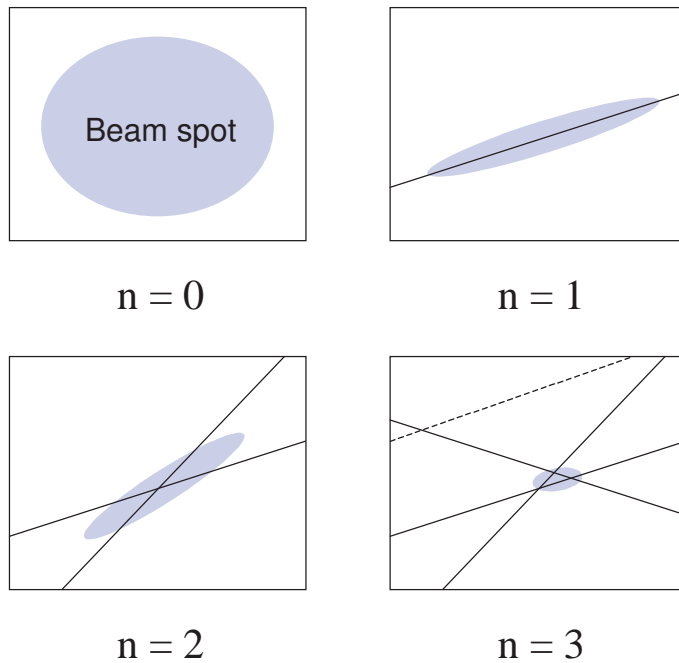


Figure 58: Illustration of the iterative construction of a (primary) vertex, where n is the number of tracks used to define the vertex in each step. The shaded area indicates the covariance ellipse of the projected vertex after each step. The dashed line indicates an outlier track.

In many practical applications, the vertex is constructed by an iterative procedure as it is illustrated in fig. 58. In most cases, some a-priori knowledge about the vertex position exists, for example the shape of the beam spot, in which interactions occur in the first place. Then a first track is selected as a vertex seed, which already narrows down the covariance ellipsoid in two dimensions. When a second suitable track is added, the vertex is already closely defined in all coordinates. This provides strong rejection power against off-vertex particles when

⁷We neglect here the complication that the B meson is likely to decay to a final state with a charmed particle which again has a non-negligible lifetime.

adding more tracks.

As in the track pattern recognition case, the danger lies in the dependence on the starting point. It is therefore necessary to use iterative criteria which ensure that the track forming the vertex seed is well chosen, and even then it must be possible to scrutinize the track ensemble of a vertex, to remove tracks that have turned out to be off the mark, and to reconnect tracks that had been discarded at an earlier stage of the construction. The vertex algorithm used in the ZEUS experiment [87], which internally uses the fitting methods of [88] may serve as an example: it uses the proton beam line as a soft constraint, and then produces a set of all track pairs that would be compatible with a common vertex together with the beam line constraint within a suitable χ^2 margin. The track pairs are then ordered according to their degree of compatibility with other track pairs, defined by the criterion above. The track pair of highest compatibility forms then the first vertex seed to be used, though also other track pairs of high compatibility level are tried, and in the end the best set is chosen based on a criterion of number of tracks and total χ^2 . Other approaches start by connecting all tracks to a diffuse master vertex, which is then successively split into vertices of smaller multiplicities and isolated tracks. A systematic investigation of different methods for vertex reconstruction in the context of the CMS experiment can be found in [89].

An entirely different approach is pursued in the topological vertex finding algorithm [90] developed for the vertex detector of the SLD experiment [17]. This method assigns a *Gaussian tube* around each track extrapolation to indicate the likelihood of an assigned vertex on a single track basis. The Gaussian tubes of all tracks are then combined to find points with maximum probability of a vertex. This method resembles the *Fuzzy Radon Transform* for tracks discussed in section 3.2. The search for maxima is then performed by sophisticated clustering algorithms. A particularly intriguing feature is the efficient resolution of heavy flavour cascade decays.

Direct vertex search by Hough transform is possible in cases where the vertex location is already strongly constrained in some coordinates, for example through the shape of a wire target [91].

6.2 Vertex Fitting

The least-squares principle can also be readily applied for vertex fitting [92, 93, 94]. The parameters of the tracks $\vec{p}_1 \dots \vec{p}_n$ at a given reference surface plus the a-priori knowledge of the vertex are the input, and the calculated vertex position together with the *reduced track parameters* of each particle, which contain only directional and momentum information at the common vertex, are the output. A general property of vertex fitting is the fact that, unlike track fitting, the fit is always non-linear, since even with straight-line tracks the extrapolation to the vertex introduces a coupling between positional and directional parameters.

As noted earlier, already vertex pattern recognition requires incremental, progressive fitting, with tracks added or removed one by one. It is therefore not surprising that also for vertex fitting, the Kalman filter is in many cases the method of choice [95]. In the vertex fitting case, the *transport* becomes trivial, and also process noise does not have an equivalent. The filter step adds another track to the vertex and updates the vertex position as well as the reduced track parameters. It is very easy to remove an already filtered track from the vertex candidate, since in the filter equations, the inverse covariance matrix of the track acts as the *weight* of the track information, and setting its sign to negative will subtract the track from the vertex fit. We prefer not to display the Kalman filter equations for vertex fitting here explicitly, but refer to the literature [27].

6.3 Kinematical Constraints

Pattern recognition deals with merging of measured information with a-priori knowledge. For example, in track pattern recognition the track model enhances the measurement power of each individual hit, while vertex assignment improves the spatial information of each associated track. In similar fashion, a-priori knowledge can be used in many cases in the further reconstruction of the event. A typical example is the beam energy constraint: in e^+e^- b-physics experiments which operate at the $\Upsilon(4S)$ energy, as BaBar, BELLE, CLEO and the earlier ARGUS, the B mesons are produced in an exclusive decay of the $\Upsilon(4S)$ resonance, and the energy of the B mesons is precisely the beam energy, which is known to a much better precision than the B meson energy reconstructed from its measured decay particles. Imposing the beam energy constraint improves then also the resolution of the B candidate mass; this method has been a vital tool in the investigation of exclusive B decays (see for example [96]).

Also masses of intermediate particles in a decay chain, for example $B^0 \rightarrow D^{*+}\pi^+\pi^-\pi^-$, $D^{*+} \rightarrow D^0\pi^+$, $D^0 \rightarrow K^-\pi^+$ can be used to imply kinematical constraints. In this case, the D^0 is a rather stable particle whose width is too small to resolve by direct kinematical reconstruction in a spectrometer. Therefore, the established knowledge of the D^0 mass [97] can be imposed as a kinematical constraint. For example, if $\vec{\alpha}$ denotes the reconstructed parameters of the K^- and π^+ particles and V_α their covariance matrix, the reconstructed D^0 mass will be a function $M(\alpha)$ of these parameters, and introduction of a Lagrange multiplier μ leads to the expression

$$X^2 = (\vec{\alpha}_c - \vec{\alpha})^T V_\alpha^{-1} (\vec{\alpha}_c - \vec{\alpha}) + 2\mu(M(\vec{\alpha}_c) - m_{D^0}) \quad (55)$$

which has to be minimized with respect to the constrained parameters $\vec{\alpha}_c$. If the daughter particles form a secondary vertex, its parameters can be optimized as well. The D^0 mass constraint leads in general to a considerable improvement of the D^* mass peak, which becomes much narrower than the experimental resolution. In comparison to the popular *mass difference method*, which benefits from

the correlation in the errors of the reconstructed D and D^* masses, this approach has the advantage that the result can be used in turn to reconstruct more complex decay chains of angular excitations in the D systems, or of B hadrons. In a next step of B reconstruction, even the tabulated D^* mass could be imposed as another independent constraint.

7 Concluding Remarks

The variety of pattern recognition tasks in particle physics tracking detectors has lead to a multitude of different approaches. Several of the global methods, as template matching or Hough transform/histogramming play an unchallenged rôle in special applications, while Hopfield networks and deformable templates frequently appear to be either limited to favourable scenarios (e.g. with 3D measurements and moderate occupancy), or need an excellent initialization or combination with a track following algorithm to become applicable at production scale. In the case of elastic arms, also the choice of an efficient minimization technique is essential. Local methods of pattern recognition are still going strong, with the Kalman filter as the mathematical backbone, and accompanied by subtle arbitration techniques they can cope well even with high track densities and sizable amounts of material in the tracking area. The new generation of high energy hadron colliders, in particular the LHC with huge track densities in piled-up events will become an important benchmark for algorithm performance. It can be expected that sophisticated combination of both global and local approaches in different passes of the procedure, matched to the particular layout of each experiment, will become a promising path to achieving the best performance.

The increasing abundance of material in radiation hard detectors poses also additional challenges to track fitting. While the correction of multiple scattering with the Kalman filter has become the accepted general standard, Molière scattering tails require a careful interpretation of the results. Electron energy reconstruction with sizable radiative energy loss is a major challenge and requires very careful treatment, and becomes a rewarding subject for robust methods beyond least square estimation. Also vertex pattern recognition can be expected to receive increasing attention in very complex event topologies at LHC, where reliable tagging of heavy flavour is a crucial prerequisite to scientific discovery.

Acknowledgement

It is a pleasure to thank E. Lohrmann for his valuable comments on the manuscript.

References

- [1] H. Grote, *Review of Pattern Recognition in High Energy Physics*, Reports on Progress in Physics 50 (1987) 473-500.
- [2] H. Albrecht et al., *Search for Rare B Decays*, Phys. Lett. B 353 (1995) 554-562.
- [3] ATLAS Collaboration, *ATLAS Inner Detector Technical Design Report Vol.I*, CERN/LHCC/97-16, CERN (1997).
- [4] C. Grupen, *Particle Detectors*, Cambridge Monographs on Particle Physics, 1996.
- [5] K. Kleinknecht, *Detectors for Particle Radiation*, Cambridge University Press (1999).
- [6] D. Green, *The Physics of Particle Detectors*, Cambridge University Press (2000)
- [7] R.L. Gluckstern, *Uncertainties in Track Momentum and Direction, Due to Multiple Scattering and Measurement Errors*, Nucl. Instr. and Meth. 24 (1963) 381-389.
- [8] R. Carlin et al. (ZEUS Collaboration), *The ZEUS Microvertex Detector*, Nucl. Instr. and Meth. A511 (2003) 23-37.
- [9] M. Danilov et al., *The ARGUS Drift Chamber*, Nucl. Instr. and Meth. 217 (1983) 153-159.
- [10] G. Sciolla et al. (BaBar Collaboration), *The BaBar Drift Chamber*, Nucl. Instr. and Meth. A419 (1998) 310-314.
- [11] F. Bruyant, J.M. Lesceux and R.J. Plano, *The Butterfly Drift Chamber Geometry: an Optimal Four-Plane Drift Chamber for Use in a High Track Multiplicity Environment*, Nucl. Instr. and Meth. 176 (1980) 409.
- [12] O. Kind et al., *A ROOT-Based Client-Server Event Display for the ZEUS Experiment*, Proc. Computing in High Energy Physics Conference, La Jolla 2003, Preprint hep-ex/0305095.
- [13] P. Krizan et al., *HERA-B, an Experiment to Study CP Violation at The HERA Proton Ring Using an Internal Target*, Nucl.Instr. and Meth. A351 (1994) 111-131.
- [14] E. Hartouni et al., *HERA-B: an Experiment to Study CP Violation in The B System Using an Internal Target at The HERA Proton Ring. Design Report*, DESY-PRC-95-01 (1995).

- [15] R. Mankel, *The HERA-B Experiment: Overview And Concepts*, Proc. International Conference on High-Energy Physics (ICHEP 98), Vancouver 1998 (Canada), Vol. 2, 1513–1518.
- [16] H. Wieman et al., *STAR TPC at RHIC*, IEEE Trans. Nucl. Sci. NS-44 (1997) 671–678.
 J.H. Thomas, *A TPC for Measuring High Multiplicity Events at RHIC*, Nucl. Instr. and Meth. A478 (2002) 166–169.
 M. Anderson et al., *The STAR Time Projection Chamber: a Unique Tool for Studying High Multiplicity Events at RHIC*, Nucl. Instr. and Meth. A499 (2003) 659–678.
- [17] F.E. Taylor et al, *Design and Performance of a 307 Million Pixel CCD Vertex Detector*, Proc. 28th International Conference on High Energy Physics (ICHEP 96), vol. 2* 1739–1742, Warsaw, (1996).
- [18] S. Brandt, *Datenanalyse* (in German), Bibliographisches Institut Mannheim 1992.
- [19] V. Blobel and E. Lohrmann, *Statistische und Numerische Methoden der Datenanalyse* (in German), Teubner (1998)
- [20] P. Bevington and D. Robertson, *Data Reduction and Error Analysis for the Physical Sciences*, McGraw/Hill 1992.
- [21] W.T. Eadie et al., *Statistical Methods in Experimental Physics*, North-Holland 1971.
- [22] A.G. Frodesen, O. Skjeggstad and H. Tofte, *Probability and Statistics in Particle Physics*, Universitetsforlaget 1979.
- [23] R.E. Kalman, *A New Approach to Linear Filtering and Prediction Problems*, Transactions of ASME Journ. Basic Engineering 82 (1960) 35–45.
 R.E. Kalman and R.S. Bucy, *New Results in Linear Filtering Prediction Theory*, Transactions of ASME Journ. Basic Engineering 83 (1961) 95–108.
- [24] P. Billoir, *Track Fitting with Multiple Scattering: a New Method*, Nucl. Instr. and Meth. 225 (1984) 352–366.
- [25] R. Frühwirth, *Application of Kalman Filtering*, Nucl. Instr. and Meth. A262 (1987) 444–450.
- [26] P. Billoir, S. Qian, *Simultaneous Pattern Recognition and Track Fitting by the Kalman Filtering Method*, Nucl. Instr. and Meth. A294 (1990) 219–228.

- [27] R.H. Böck, H. Grote, D. Notz and M. Regler, *Data Analysis Techniques for High-Energy Physics Experiments*, Cambridge Univ. Press (1990). 2nd edition (with R. Frühwirth) (2000)
- [28] D.N. Brown, E.A. Charles and D.A. Roberts, *The BaBar Track Fitting Algorithm*, Proc. Computing in High Energy Physics Conference, Padova (2000).
- [29] R. Mankel and A. Spiridonov, *Compatibility Analysis*, HERA-B Internal Note 99-111 (1999).
- [30] H.D. Schulz and H.J. Stuckenberg, Proc. Topical Conference on the Application of Microprocessors in High Energy Physics Experiments, CERN 81-07 (1981).
- [31] N. Koch et al., *The ARGUS Vertex Trigger*, Nucl. Instr. and Meth. A373 (1996) 387-405.
- [32] S. Seidel et al. (ARGUS Collaboration), *The ARGUS Micro-Vertex Drift Chamber*, Proc. APS Conference Particles and Fields, Vol. 2, 1158-1163, Vancouver (1991).
- [33] M. Dell'Orso and L. Ristori, *A Highly Parallel Algorithm for Track finding*, Nucl. Instr. and Meth. A287 (1990) 436-438.
- [34] P. Battaiotto et al., *The Tree-Search Processor for Real-Time Track Pattern Recognition*, Nucl. Instr. and Meth. A287 (1990) 431-435.
- [35] K. Ackerstaff et al., *The HERMES Spectrometer*, Nucl. Instr. and Meth. A417 (1998) 230-265.
- [36] J. Blom et al., *A Fuzzy Radon Transform for Track Recognition*, Proc. Computing in High Energy Physics Conference, San Francisco (1994).
- [37] M. Gyulassy and M. Harlander, *Elastic Tracking and Neural Network Algorithms for Complex Pattern Recognition*, Comp. Phys. Comm. 66 (1991) 31-46.
- [38] A. Antonov (Moscow Engineering and Physics Institute), private communication.
- [39] P.V.C. Hough, *Machine Analysis of Bubble Chamber Pictures*, Int. Conf. on High Energy Accelerators and Instrumentation, 554-556, CERN, 1959.
- [40] M. Ohlsson, C. Peterson and A.L. Yuille, *Track Finding with Deformable Templates: The Elastic Arms Approach*, Comput. Phys. Commun. 71 (1992) 77-98.

- [41] C. Borgmeier, *Global Pattern Recognition in the HERA-B Tracking System* (in German), Diploma thesis, Humboldt University Berlin (1996).
- [42] T. Schober, *Investigation of Hough Transforms as Global Approaches to Pattern Recognition in the HERA-B Main Tracking System* (in German), Diploma thesis, Humboldt University Berlin (1996).
- [43] J.A.Anderson and E. Rosenfeld, *Neurocomputing: Foundations of Research*, MIT Press, Cambridge (1988).
- [44] J.J. Hopfield, *Neural Networks and Physical Systems with Emergent Collective Computational Abilities*, Proc. National Academy of Science, USA, 79 (1982) 2554-2558. Reprinted in [43].
- [45] Y. Shrivastava, S. Dasgupta and S.M. Reddy, *Guaranteed Convergence in a Class of Hopfield Networks*, IEEE Transactions on Neural Networks Vol.3, No.6 (1992) 951-961. Reprinted in [43].
- [46] B. Denby, *Neural Networks and Cellular Automata in Experimental High energy Physics*, Comput. Phys. Commun. 49 (1988) 429-448.
- [47] C. Peterson, *Track Finding with Neural Networks*, Nucl. Instr. and Meth. A279 (1989) 537-549.
- [48] G. Stimpff-Abele and L. Garrido, *Fast Track Finding With Neural Networks*, Comput. Phys. Commun. 64 (1991) 46-56.
- [49] R. Mankel, *Pattern Recognition Algorithms For B Meson Reconstruction in Hadronic Collisions*, Proc. Computing in High Energy Physics Conference, Berlin (1997). URL <http://www.ifh.de/CHEP97/paper/183.ps>
- [50] I. Abt et al., *Cellular Automaton And Kalman Filter Based Track Search in The HERA-B Pattern Tracker*, Nucl. Instr. and Meth. A490 (2002) 546-558.
- [51] I. Abt et al., *CATS: a Cellular Automaton for Tracking in Silicon for The HERA-B Vertex Detector*, Nucl. Instr. and Meth. A489 (2002) 389-405.
- [52] M. Ohlsson, *Extensions and Explorations of the Elastic Arms Algorithm*, Comput. Phys. Commun. 77 (1993) 19-32.
- [53] C. Peterson and B. Söderberg, *A New Method for Mapping Optimization Problems onto Neural Networks*, Int. Journal of Neural Systems 1 (1989) 3-22.
- [54] M. Lindström, *Track Reconstruction in The ATLAS Detector Using Elastic Arms*, Nucl. Instr. and Meth. A357 (1995) 129-149.

- [55] R. Blencenbecler, *Deformable Templates – Revised And Extended – With an OOP Implementation*, Comput. Phys. Commun. 81 (1994) 318-334.
- [56] A. Paus, *Pattern Recognition in the Tracking System of the HERA-B Detector With an Elastic Arms Algorithm* (in German), Diploma thesis, Humboldt University Berlin (1997).
- [57] S.E. Fahlmann, *Faster-Learning Variations on Back-Propagation: an Empirical Study*, Proc. Connectionist Models Summer School, San Mateo (1988).
- [58] A. Riedmiller and H. Braun, *A Direct Adaptive Method for Faster Back-propagation Learning: the RPROP Algorithm*, Proc. IEEE International Conference of Neural Networks, San Francisco (1993).
- [59] R. Frühwirth and A. Strandlie, *Track Fitting with Ambiguities and Noise: a Study of Elastic Tracking and Nonlinear Filters*, Comp. Phys. Comm. 120 (1999) 197-214.
- [60] A. Strandlie and R. Frühwirth, *Adaptive Multitrack Fitting*, Comp. Phys. Comm. 133 (2000) 34-42.
- [61] R. Mankel, *A Concurrent Track Evolution Algorithm for Pattern Recognition in the HERA-B Main Tracking System*, Nucl. Instr. and Meth. A395 (1997) 169-184.
- [62] H. Albrecht, private communication.
- [63] A. Khanov et al., *Tracking in CMS: Software Framework And Tracker Performance*, Nucl. Instr. and Meth. A478 (2002) 460-464.
- [64] M.M. Angarano et al. (CMS Tracker Collaboration), *The Silicon Strip Tracker for CMS*, Nucl. Instr. and Meth. A501 (2003) 93-99.
- [65] R. Mankel and A. Spiridonov, *The Concurrent Track Evolution Algorithm: Extension for Track Finding in The Inhomogeneous Magnetic Field of The HERA-B Spectrometer*, Nucl. Instr. and Meth. A426 (1999) 268-282.
- [66] M. Regler, R. Frühwirth and W. Mitaroff, *Filter Methods in Track And Vertex Reconstruction*, Int. Journ. Mod. Phys. C7 (1996) 521-542.
- [67] R. Mankel, *Online Track Reconstruction for HERA-B*, Nucl. Instr. and Meth. A384 (1996) 201-206.
- [68] W.H. Press et al., *Numerical Recipes in C: the Art of Scientific Computing*, 2nd edition, Cambridge University Press, 1993.

- [69] T. Oest, *Particle Tracing Through The HERA-B Magnetic Field*, HERA-B Internal Note 97-165 (1997).
- [70] H.A. Bethe, *Molière's Theory of Multiple Scattering*, Phys. Rev. 89 (1953) 1256-1266.
- [71] V.L. Highland, *Some Practical Remarks on Multiple Scattering*, Nucl. Instr. and Meth. 129 (1975) 497-499.
Erratum Nucl. Instr. and Meth. 161 (1979) 171.
- [72] G.R. Lynch and O.L. Dahl, *Approximations for Multiple Coulomb Scattering*, Nucl. Instr. and Meth. B58 (1991) 6-10.
- [73] G. Lutz, *Optimum Track Fitting in The Presence of Multiple Scattering*, Nucl. Instr. and Meth. A273 (1988) 349-374.
- [74] E.J. Wolin and L.L. Ho, *Covariance Matrices for Track Fitting With The Kalman Filter*, Nucl. Instr. and Meth. A329 (1993) 493-500.
- [75] R. Mankel, *The Object-Oriented Track Fit*, HERA-B Internal Note 98-079.
- [76] G. Baum et al. (COMPASS Collaboration), *Common Muon and Proton Apparatus for Structure and Spectroscopy (Proposal)*, CERN/SPSLC 96-14.
- [77] S. Amato et al. (LHCb Collaboration), *LHCB Technical Proposal*, CERN-LHCC-98-4, CERN-LHCC-P-4 (1998).
- [78] V. Papavassiliou et al. (BTeV Collaboration), *BTeV: A Proposal for a New B Physics Experiment at the Fermilab Tevatron Collider La Thuile 2000*, Results And Perspectives in Particle Physics, 843-864 (2000)
- [79] D. Stampfer, M. Regler and R. Frühwirth, *Track Fitting With Energy Loss*, Comput. Phys. Commun. 79 (1994) 157-164.
- [80] H.A. Bethe and W. Heitler, Proc. Roy. Soc. A146 (1934) 83.
- [81] G. Kitagawa, *Non-Gaussian Seasonal Adjustment*, Comp. and Math. Appl. 18 (1989) 503-514.
- [82] G. Kitagawa, *The Two-Filter Formula for Smoothing And an Implementation of The Gaussian-Sum Smoother*, Annals Inst. Statist. Math. 46 (1994) 605-623.
- [83] R. Frühwirth, *Track Fitting With Long-Tailed Noise: a Bayesian Approach*, Comput. Phys. Comm. 85 (1995) 189-199.

- [84] R. Frühwirth, *Track Fitting With Non-Gaussian Noise*, Comput. Phys. Comm. 100 (1997) 1-16.
- [85] R. Frühwirth and M. Regler, *On the Quantitative Modelling of Tails and Core of Multiple Scattering by Gaussian Mixtures*, Nucl. Instr. and Meth. A456 (2001) 369.
- [86] R. Frühwirth and S. Frühwirth-Schnatter, *On the Treatment of Energy Loss in Track Fitting*, Comput. Phys. Comm. 110 (1998) 80-86.
- [87] G.F. Hartner, *VCTRAK Briefing: Program and Math*, ZEUS Internal Note 98-058, 1998.
- [88] P. Billoir and S. Qian, *Fast Vertex Fitting with Local Parametrization of Tracks*, Nucl.Instr. and Meth. A311 (1992) 139-150.
- [89] R. Frühwirth et al., *New Vertex Reconstruction Algorithms for CMS*, Proc. Computing in High Energy Physics Conference, La Jolla 2003, Preprint physics/0306012.
- [90] D.J. Jackson, *A Topological Vertex Reconstruction Algorithm for Hadronic Jets*, Nucl. Instr. and Meth. A388 (1997) 247-253.
- [91] T. Lohse, *Vertex Reconstruction and Fitting*, HERA-B Internal Note 95-013, 1995.
- [92] D.H. Saxon, *Three-Dimensional Track and Vertex Fitting in Chambers with Stereo Wires*, Nucl. Instr. and Meth. A234 (1985) 258-266.
- [93] G.E. Forden and D.H. Saxon, *Improving Vertex Position Determination Using a Kinematic Fit*, Nucl. Instr. and Meth. A248 (1986) 439-450.
- [94] D.H. Saxon, *Vertex Detection and Tracking at Future Accelerators*, Hadronic J. 10 (1987) 117-139.
- [95] R. Luchsinger and C. Grab, *Vertex Reconstruction by Means of the Kalman Filter*, Comput. Phys. Comm. 76 (1993) 263-280.
- [96] H. Albrecht et al., *Exclusive Hadronic Decays of B Mesons*, Z. Phys. C48 (1990) 543-551.
- [97] Particle Data Group (K. Hagiwara et al.), *Review of Particle Physics*, Phys. Rev. D66 (2002).