

Fast neural-net based fake track rejection

Michel DeCian¹, Stephen Farry², Paul Seyfert³, Sascha Stahl⁴.

¹*Physikalisches Institut, Ruprecht-Karls-Universität Heidelberg, Heidelberg, Germany*

²*Oliver Lodge Laboratory, University of Liverpool, Liverpool, United Kingdom*

³*Sezione INFN di Milano Bicocca, Milano, Italy*

⁴*European Organization for Nuclear Research (CERN), Geneva, Switzerland*

Abstract

A neural-network based algorithm to identify fake tracks in the LHCb pattern recognition is presented. This algorithm, called ghost probability, is fast enough to fit into the CPU time budget of the software trigger farm. It allows reducing the fake rate and consequently the combinatorics of the decay reconstructions, as well as the number of tracks that need to be processed by the particle identification algorithms. As a result, it strongly contributes to the achievement of having the same reconstruction online and offline in the LHCb experiment.

1 Introduction

The LHCb detector consists of subsystems designed to perform high efficiency tracking ($> 95\%$) with an excellent momentum resolution (0.5% for $p < 20 \text{ GeV}/c$). Two Ring Imaging Cherenkov detectors provide precise particle identification. In Run II of the LHC, a new scheme for the LHCb software trigger allows splitting the triggering of the event in two stages, giving room to perform the alignment and calibration in real time. In the novel detector alignment and calibration strategy for Run II, data collected at the start of the fill are processed in a few minutes and used to update the alignment, while the calibration constants are evaluated for each run. This allows identical constants to be used in the online and offline reconstruction.

One of the challenges to achieve run the offline reconstruction in the software trigger is the limited CPU time budget of the computing farm. The reconstruction time of events depends strongly on the number of reconstructed charged particle tracks in an event in two ways. The particle identification (PID) is evaluated for every reconstructed track, and the combinatorics of reconstruction decay vertices gets more complex with more reconstructed tracks.

A key ingredient to fit the offline reconstruction into the software trigger is the reduction of the fake track rate prior to the PID and combinatorics of reconstructing particle decays in the second software trigger stage. A neural network, described in this note, is deployed to identify fake tracks, called the “ghost probability”.

21 Terminology

To avoid ambiguity, the bare term “performance” is avoided. Instead, when referring to how well good tracks are separated from fake tracks, the term “physics performance” is used since it is the figure of merit on which physics analyses depend. The term “CPU performance” is used for the amount of computing resources needed to execute the algorithm proposed in this note. As benchmark for the latter, the cycle count of callgrind [1] is used. Effects of instruction caching and data caching are assumed small, approximately confirmed by wall clock time measurements. The cycles spent in other algorithms which are only called to compute input quantities to the ghost probability are accounted to the ghost probability, most notably this comprises the interpolation of tracks through active detector material to determine which channels should have a hit from the track – algorithms like the track fit, which would be executed anyways, are not accounted to the ghost probability.

The term “ghost probability” is used for both, the entire algorithm computing whether a track is considered a fake track or a real track, including the neural network, and for the numeric response of that algorithm. When the ghost probability is referred to as a selection requirement, the nominal working point of 0.4 is implied, corresponding to a fake track retention of 40%.

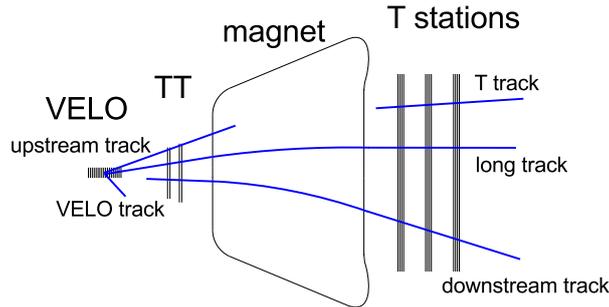


Figure 1: Illustration of the tracking system of LHCb, starting from the VELO around the collision point on the left, particles pass the TT, are deflected in the magnetic field of the dipole magnet and then detected in the T station detector (IT and OT). Different track types are reconstructed by different track finding algorithms. [2]

39 1.1 Track reconstruction

40 Owing to the design of the LHCb detector, which consists of tracking detectors mainly
 41 outside the magnetic field, charged particle tracks are in approximation straight line
 42 segments in the upstream part (VELO and TT) and in the downstream part (T stations).
 43 Figure 1 shows an overview of the different track types defined in the LHCb reconstruction:
 44 VELO tracks, which have hits in the VELO only; upstream tracks, which have hits in
 45 the two upstream trackers; T tracks, which have hits in the T stations only; downstream
 46 tracks, which have hits in TT and the T stations; and long tracks, which have hits in the
 47 VELO and the T stations. The latter tracks can additionally have hits in TT.

48 If a particle is reconstructed more than once, as different track types, only the track
 49 best suited for analysis purposes is kept. Hereby, long tracks are preferred over any other
 50 track type, upstream tracks are preferred over VELO tracks, and downstream tracks are
 51 preferred over T tracks.

52 Most analyses use long tracks because they provide the best momentum and spatial
 53 resolution among all track types. Unless otherwise stated, track reconstruction at LHCb
 54 refers to the reconstruction of long tracks. In a typical signal triggered event, around 60
 55 long tracks are reconstructed. Other track types, such as downstream tracks, are used
 56 for the reconstruction of decay products of long-lived particles such as K_s^0 mesons, or for
 57 internal alignment of the tracking detectors.

58 Tracks are fit with a Kalman filter fit. In addition to a global fit χ^2 , separate
 59 contributions to the χ^2 from the downstream detectors (IT and OT), χ_D^2 , and from the
 60 upstream detectors (VELO and TT), χ_U^2 are computed. A large number of fake tracks
 61 results from wrong combinations of well reconstructed track segments in the upstream and
 62 downstream regions. These usually have good χ_D^2 and χ_U^2 but the additional contribution
 63 from matching the two segments, $\chi_M^2 = \chi^2 - \chi_D^2 - \chi_U^2$, is large for these “matching” fakes.

64 The Kalman fit has an outlier removal to account for individual detector hits which are
 65 not due to the reconstructed particle track. Beyond that, a special treatment for Outer

66 Tracker hits is in place. The readout electronics is designed to select only a single hit in
67 each channel per bunch crossing; if two charged particles pass the same straw, a drift time
68 measurement will only be provided for one of them. To describe tracks in high occupancy
69 OT modules, the drift time measurement can be ignored and only the information that
70 a track went somewhere through the straw is used. This is decided for each straw–track
71 combination individually if the hit residual is too large, similar to a standard outlier
72 removal. This drift time suppression ensures that the track fit χ^2 is not biased to larger
73 values for tracks in high multiplicity events, for tracks in the OT with respect to tracks in
74 the IT, or for tracks in high occupancy modules, which are those closer to the beam axis.

75 1.2 Previous works

76 An earlier version of the work presented here, referred to as old ghost probability, was
77 already used in analyses of Run I data. The neural network was evaluated in the offline
78 reconstruction to distinguish fake tracks from real particles’ tracks [3] (used e.g. in [4]).
79 The network was trained on all reconstructed tracks in simulated events with at least one
80 $b\bar{b}$ pair produced in the pp collision.

81 The 22 input variables to the old ghost probability are the track fit χ^2 , and the individual
82 contributions $\chi_D^2, \chi_U^2, \chi_M^2$ and the corresponding degrees of freedom; the numbers of hits
83 on the track in each tracking detector; the reconstructed track p_T and pseudorapidity; the
84 difference in the number of observed hits on a track and the “expected hits”, calculated
85 interpolating the track through the detector and counting how many active strips/straws
86 the track passes through; and finally the occupancies of all tracking detectors.

87 There are separate networks for each track type, where input variables are removed if
88 they are not defined for that track type (e.g. VELO hits for downstream tracks).

89 1.3 Network architecture tuning

90 As framework for the neural network, the TMVA package [5] is chosen since it is

- 91 • equipped with a root file interface for the training, which is the common data file
92 format in LHCb software,
- 93 • commonly known in LHCb (ensuring future maintainability),
- 94 • able to provide code generation for the trained network such that the network can
95 be integrated into any C++ code without creating dependency on external libraries.

96 Mathematically, the shallow neural network that is implemented a composed function

$$\mathbb{R}^n \xrightarrow{\text{linear} + \text{const}} \mathbb{R}^m \xrightarrow{\text{element wise non-linear}} \mathbb{R}^m \xrightarrow{\text{linear} + \text{const}} \mathbb{R} \xrightarrow{\text{non-linear}} \mathbb{R}.$$

97 The entries of the matrices for the linear mappings ($(n+1) \times m$ for the first and $(m+1)$
98 for the second) are subject to optimisation, where n is the number of input variables and

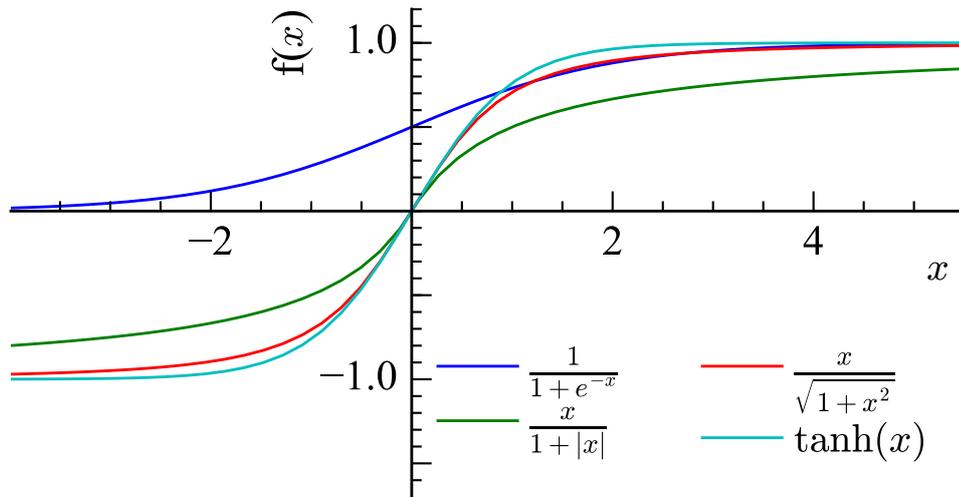


Figure 2: Functional shape of sigmoid functions.

99 m chosen as $n + 5$. The non-linearities, so-called activation functions, are fixed real valued
 100 functions.

101 At the time of development, the $\tanh(x)$ function was a commonly used activation
 102 function in TMVA, while known as a computationally expensive function to be optimised
 103 for the LHCb pattern recognition [6]. Yet it is not the only possible sigmoid function [7]
 104 and consequently custom activation functions have been added to TMVA [8].

105 Of the tested functions $\frac{x}{\sqrt{1+x^2}}$ is the fastest to compute, while no significant physics
 106 performance difference is expected given the similar functional shape, see Fig. 2. Indeed,
 107 no performance difference is observed in Fig. 3. Therefore, it is chosen as activation
 108 function.

109 The ghost probability is a classification problem, and thus cross entropy [9] is chosen
 110 as loss function in the network training. With respect to the run I implementation of the
 111 ghost probability, this contributes to the physics performance improvement. The activation
 112 function of the output layer is $\frac{1}{1+e^{-x}}$ in the training. In the application, a custom output
 113 calibration is applied instead, as described in Sect. 1.5).

114 1.4 Variable selection

115 To allow for enough development time for testing and evaluation, the selection of input
 116 variables is mostly unchanged from Run I with two exceptions. The track interpolation to
 117 determine the number of expected hits is removed to reduce the CPU usage of the ghost
 118 probability by a factor 10. The number of track candidates competing for shared hits in
 119 the pattern recognition is added as input variable.

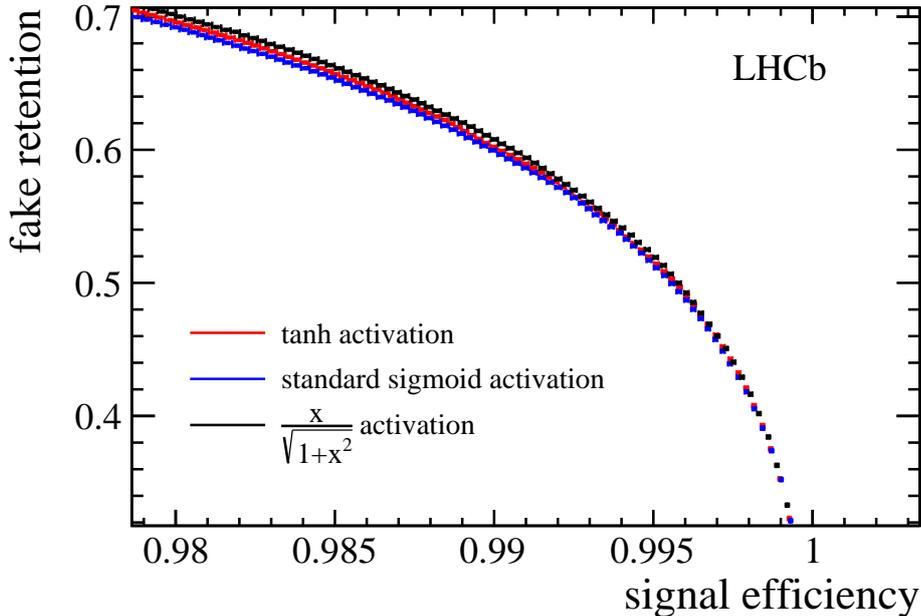


Figure 3: ROC curves for fake track discriminating neural networks, using different activation functions. A very small physics performance improvement is observed when changing from the TMVA standard functions to $\frac{x}{\sqrt{1+x^2}}$.

120 1.5 Output transformation

121 To ease the usage of the ghost probability, a transformation of the network response is
 122 applied. A probability integral transform – also referred to as “flattening” or “rarity
 123 transformation” – is obtained as a linear spline fit to the cumulative network response for
 124 fake tracks in simulated events. The discriminating behaviour of any classifier is invariant
 125 under monotonous transformations and so is the physics performance under the probability
 126 integral transform. Motivations for this transformation are primarily to give a physical
 127 interpretation to the response: rejecting tracks with a ghost probability of larger than $x\%$
 128 will retain $x\%$ of all fake tracks.

129 In addition any update of the ghost probability training will have the same behaviour
 130 and thus the optimal working points of algorithms downstream of the ghost probability
 131 algorithm will remain unchanged at leading order.

132 1.6 Category classifiers

133 Fake tracks produced by different pattern recognition algorithms might have different
 134 track properties. It might therefore be beneficial to train separate neural networks for
 135 the two main track reconstruction algorithms at LHCb. On simulated events, the physics
 136 performance of two separate networks does not differ from the physics performance of a
 137 single network. Similarly, different networks for different T station regions have been tested

138 (one for tracks in the OT, IT, and the overlap region), without significant performance
139 gain. We suspect reasons are twofold. Firstly, the network already knows which T station
140 tracker a track went through due to the hit counts in the individual subdetectors. Secondly,
141 for the different algorithms, we suspect the common track fit is more indicative for whether
142 a track is a fake or not, than the pattern recognition strategy.

143 Consequently, a single network for all pattern recognition algorithms in the entire
144 detector is deployed.

145 1.7 Training sample

146 The LHCb track reconstruction needs to be able to handle a wide range of LHC running
147 conditions. At the time of preparing for data taking in 2015 it was not clear whether the
148 LHC would operate at 25 ns bunch spacing or 50 ns bunch spacing. Simulations to prepare
149 the track reconstruction were prepared for these scenarios:

- 150 • 25 ns bunch spacing, $\nu = 1.6$
- 151 • 25 ns bunch spacing, $\nu = 1.9$
- 152 • 50 ns bunch spacing, $\nu = 1.6$
- 153 • 50 ns bunch spacing, $\nu = 2.7$

154 where ν is the average number of pp interactions per bunch crossing.

155 The scenarios differ, for what concerns the track reconstruction, significantly in detector
156 occupancy and spillover in the Outer Tracker. That may lead to different behaviours
157 of fake track reconstruction and require different network trainings for different running
158 conditions. The necessity for having different network trainings is assessed by training
159 networks for each of the running conditions, with all other training parameters fixed,
160 and evaluating the networks on one of the samples and their discriminating powers are
161 compared. Figure 4 shows the ROC curves for the 25 ns sample at low pile-up. The
162 discrimination powers of the four networks do not largely differ and thus for simplicity only
163 a single network (trained at the favoured scenario of 25 ns bunch spacing at low pile-up) is
164 deployed.

165 2 Validation

166 The data taking strategy of LHCb in Run II involves the application of the same track
167 reconstruction in the software trigger as in the offline data processing. This goal can
168 only be achieved within the time constraints of the software trigger by applying the ghost
169 probability in the trigger. This ghost probability reduces the rate of fake tracks entering
170 the particle identification and combinatorics of decay reconstructions and thereby saves
171 more time than the computation of the ghost probability.

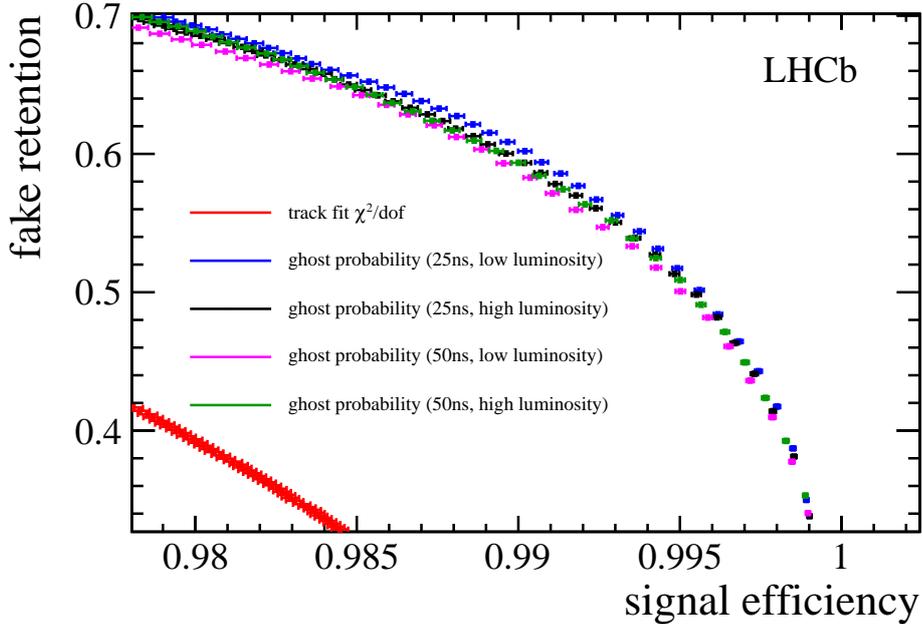


Figure 4: ROC curves for fake track discriminating neural networks (blue, black, magenta, and green), trained for different LHC running conditions, evaluated for 25 ns bunch spacing, $\nu = 1.6$. The red points are the ROC curve for the track fit reduced χ^2 , which performs significantly worse than any of the neural networks.

172 It must therefore be ensured that the full physics program of LHCb can be done with
 173 tracks passing the ghost probability, and that there is no corner of phase space or particle
 174 species, which is rejected by the ghost probability.

175 The computation of the track fit χ^2 was last revised in 2015 between data taking
 176 at 50 ns and 25 ns bunch spacing, [10]. For conclusive validations, tracks in Run I data
 177 and from 2015 data with 50 ns bunch spacing are refitted before computing the ghost
 178 probability.

179 2.1 High momentum tracks

180 Due to their low cross section, tracks in the momentum range of $Z \rightarrow \mu\mu$ events are absent
 181 in the training data. This could lead to a low selection efficiency for very high momentum
 182 tracks.

183 In the early measurement period in 2015 at a bunch spacing of 50 ns, the nominal 2015
 184 pattern recognition was used without application of the ghost probability. Refitting the
 185 candidate tracks from $Z \rightarrow \mu\mu$ decays in that period allows to assess the performance of
 186 the ghost probability for very high momentum tracks. The measured efficiency in Fig. 5
 187 shows that the absence of very high momentum tracks does not lead to a low efficiency.

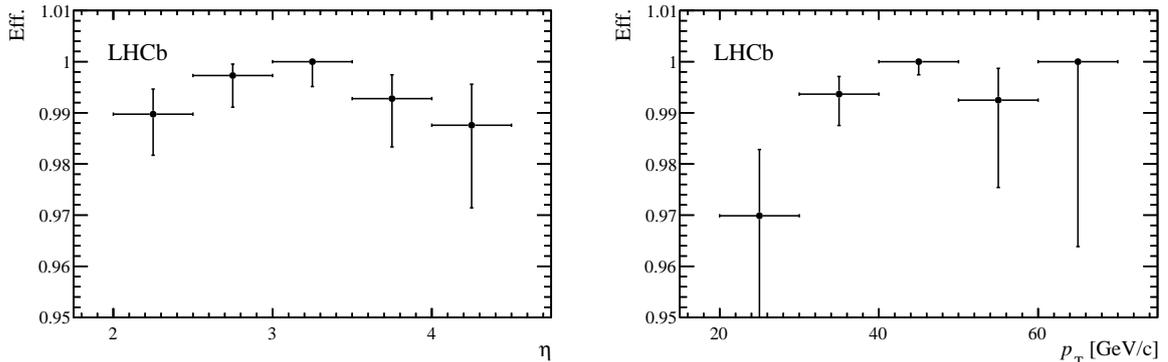


Figure 5: Efficiency for $Z \rightarrow \mu\mu$ tracks to pass the ghost probability, in data from 2015 at a bunch spacing of 50 ns.

188 2.2 Electron reconstruction

189 Electrons are more challenging to reconstruct than the standard candles ($Z \rightarrow \mu\mu$ or
 190 $D \rightarrow K\pi$). At the same time, it can be expected that the response of the ghost probability
 191 for electrons differs from that for other particles as electrons undergo more multiple
 192 scattering.

193 It is assumed that the reconstruction of converted photons as electron pair is the most
 194 vulnerable part for the following reasons. The photon conversion can happen “late” in the
 195 VELO leaving only few hits. In addition, the e^+e^- pair has a small opening angle which
 196 could lead to hit ambiguities in the VELO pattern reconstruction. It should be noted that
 197 analyses of channels like $B_s^0 \rightarrow K^*\gamma$ are anyhow so-called rare decays which immediately
 198 suffer from efficiency loss.

199 The “early data” of 2015 does not correspond to enough integrated luminosity to
 200 obtain a satisfying estimation of the consequences of a cut on the ghost probability on
 201 converted photons. For this reason, the tracks of $B_s^0 \rightarrow K^*\gamma$ candidates from Run I – using
 202 a simplified version of the selection presented in [11] without Bremsstrahlungs correction –
 203 are refitted using the track fit configuration as used in 2015 and the ghost probability is
 204 evaluated. The invariant mass spectrum shown in Fig. 6 shows candidates without the
 205 application of a cut on the ghost probability, those passing, and those failing; both for
 206 converted photons reconstructed as pair of downstream tracks and as pair of long tracks.
 207 To the statistical precision of this test, no signal loss is visible.

208 2.3 Validation with 25 ns data

209 A cut on the ghost probability is included in the standard track reconstruction since data
 210 taking at 25 ns bunch spacing started. To investigate the behaviour of the ghost probability
 211 in real data with 25 ns bunch spacing, events are re-reconstructed without a cut on the
 212 ghost probability. Under the assumption, that most K_s^0 are part of the underlying event
 213 and most triggered events containing K_s^0 would have been triggered without those K_s^0 , K_s^0

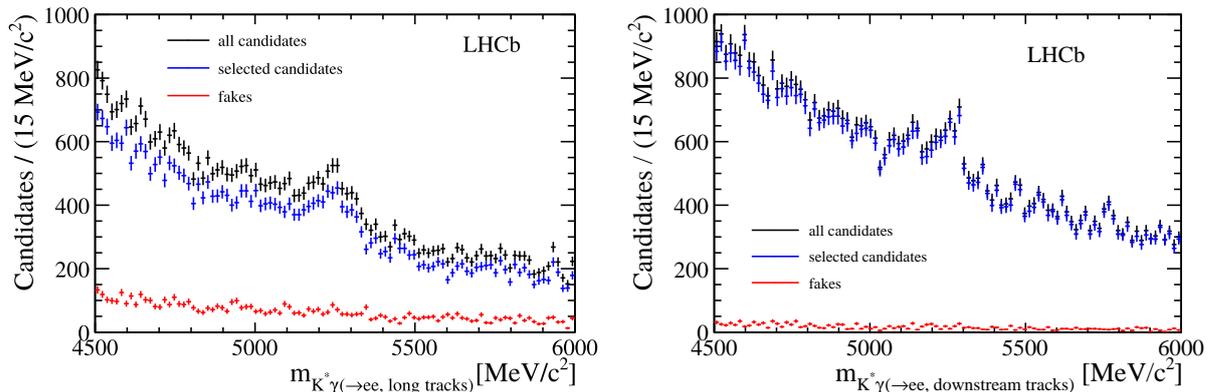


Figure 6: $B_s^0 \rightarrow K^* \gamma$ candidates, where the photon is reconstructed as pair of electron tracks, using long tracks (left) and using downstream tracks (right). The candidates prior to a cut on the ghost probability of the electron tracks are shown in black, those passing the ghost probability in blue, and those candidates rejected by the ghost probability in red. No signal loss is visible in the rejected candidates.

214 are used as probe of the ghost probability.

215 The invariant mass spectrum of K_s^0 candidates after re-reconstruction is shown in
 216 Fig. 7 for both long tracks and downstream tracks. In both cases, K_s^0 candidates are
 217 reconstructed from two opposite charged pions which are compatible with originating
 218 from a common vertex, which satisfy fiducial momentum requirements, and which are
 219 significantly displaced from any primary collision vertex. In either case, the background
 220 contribution is largely reduced when rejecting events where at least one of the tracks has
 221 a ghost probability of larger than 0.4. There is no signal visible in the events rejected. It
 222 is concluded that no physics signal is lost due to the application of the ghost probability
 223 within the statistical sensitivity of the test in the kinematic spectrum of the selected K_s^0 .

224 The same test with $D \rightarrow K \pi$ decays is shown in Fig. 8 on the right. Similar to the K_s^0
 225 selection, kaon and pion tracks are selected with minimal momentum requirements and
 226 are required to originate from a common, displaced vertex. Additionally, the kaon track
 227 must be identified as kaon by the RICH system. To ensure that the sample is not biased
 228 towards candidates passing the ghost probability due to the online event selection, the
 229 ghost probability spectrum is shown in Fig. 8 b), where no step from such a selection is
 230 visible at 0.4.

231 2.4 Decay time acceptance

232 The study of long lived particles (b and c hadrons) is the major part of the LHCb physics
 233 program. It must therefore be ensured that the ghost probability does not reject particles
 234 from displaced vertices at a higher probability than particles from primary collisions (which
 235 have a higher prevalence in the training).

236 To evaluate a possible decay time bias of the ghost probability, for each reconstructed

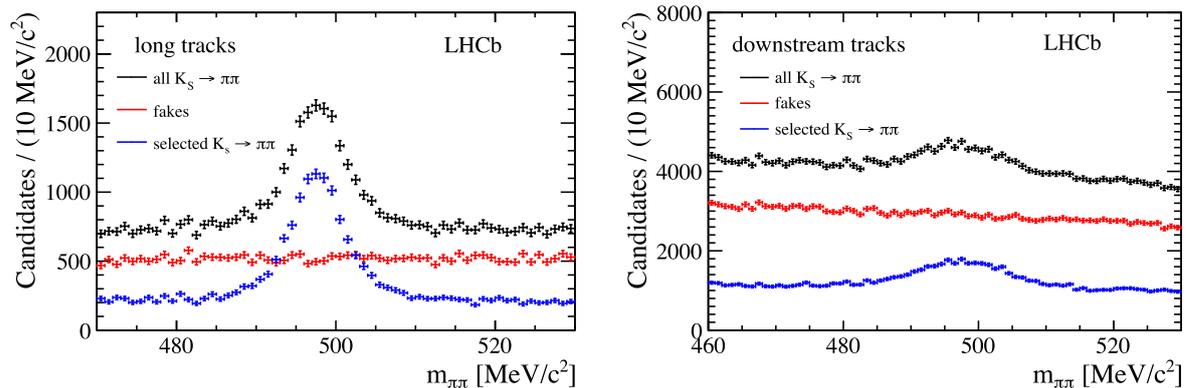


Figure 7: $K_S^0 \rightarrow \pi\pi$ invariant mass spectrum for events reconstructed without using the ghost probability in the track selection; using long tracks (left) and downstream tracks (right). Candidates for which at least one track fails the default ghost probability requirement are shown in red and do not exhibit a signal contamination. The remaining candidates are shown in blue.

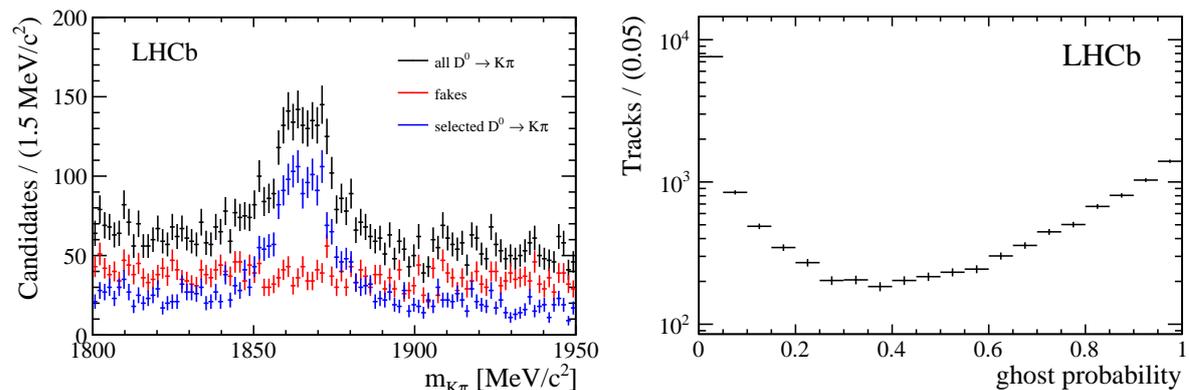


Figure 8: $D \rightarrow K\pi$ invariant mass spectrum for events reconstructed without using the ghost probability in the track selection (left). Candidates for which at least one track fails the default ghost probability requirement are shown in red and do not exhibit a signal contamination. The remaining candidates are shown in blue. The ghost probability distribution for the kaon and pion tracks is shown on the right.

237 particle in simulated events with a $b\bar{b}$ production, the average of the true decay time of their
 238 ancestor particles is determined. When rejecting tracks which fail the ghost probability
 239 criterion, the average ancestor decay time changes by $(1.5 \pm 2.0) \times 10^{-15}$ s. This is smaller
 240 than the statistical sensitivity of this test, and smaller than the systematic uncertainty to
 241 which the lifetime bias of the LHCb reconstruction is known [12].

3 Outlook

The current networks are trained for the track reconstruction for data taking in 2015 at 25 ns bunch spacing, using the latest simulations available at the time. Retraining is advisable for significant updates in the track reconstruction “upstream” of the ghost probability (i.e. the pattern recognition and the track fit). Physics performance gains can also be expected with improved machine learning techniques or event simulations.

Additional separation between “good” tracks and fake tracks could be gained by using hit expectations in active layers: at the moment only the numbers of hits in the individual subdetectors on the track are used. These could be compared with the intersections of the trajectory with active detector material such that the number of missing hits is used as input for the network [13].

The current training is purely based on simulated events, the domain adaptation approach from [14] is not applied as it currently does not lead to an improved fake track rejection. The ghost probability network is retrained using good tracks and fake tracks from simulated events and unlabelled tracks from real events with the Caffe software framework [15], which is used in [14]. In addition to the network with domain adaptation, a conventional network is trained to disentangle effects from the training algorithm (TMVA \rightarrow Caffe) and network architecture (adding a gradient reversal layer and domain classifier). This working point of the network responses is chosen to retain the same number of K_S^0 candidates as the application of the nominal ghost probability. From the invariant mass distribution in Fig. 9, it can be seen that the TMVA network and the two networks trained in Caffe yield close to identical physics performance; the data points of the network with domain adaptation are almost entirely covered by the data points of the Caffe network without domain adaptation. This does not rule out that domain adaptation can not improve the physics performance of the ghost probability in the future.

The current network relies on auto-vectorisation. The methods suggested by [16] lead to tenfold improvement of the neural network implementation [17] once using AVX intrinsic commands. This approach has not been followed up to ensure platform independence of the ghost probability.

The current activation function in the neural network is $\frac{x}{\sqrt{1+x^2}}$. The rectified linear unit $\max(0, x)$ or $\frac{x}{1+|x|}$ are expected to be even faster, as listed in Tab. 1 (from [18]).

4 Conclusion

The ghost probability is introduced as a default method of reducing the number of fake tracks in the LHCb reconstruction and is deployed for offline and online reconstruction. The reduction of fake tracks in the particle identification and combinatorics of decay reconstruction greatly reduces the demand of computing resources of the software trigger and enables LHCb to use identical reconstructions online and offline. Validations at different phase space points reveal no adverse side effects of applying the ghost probability centrally.

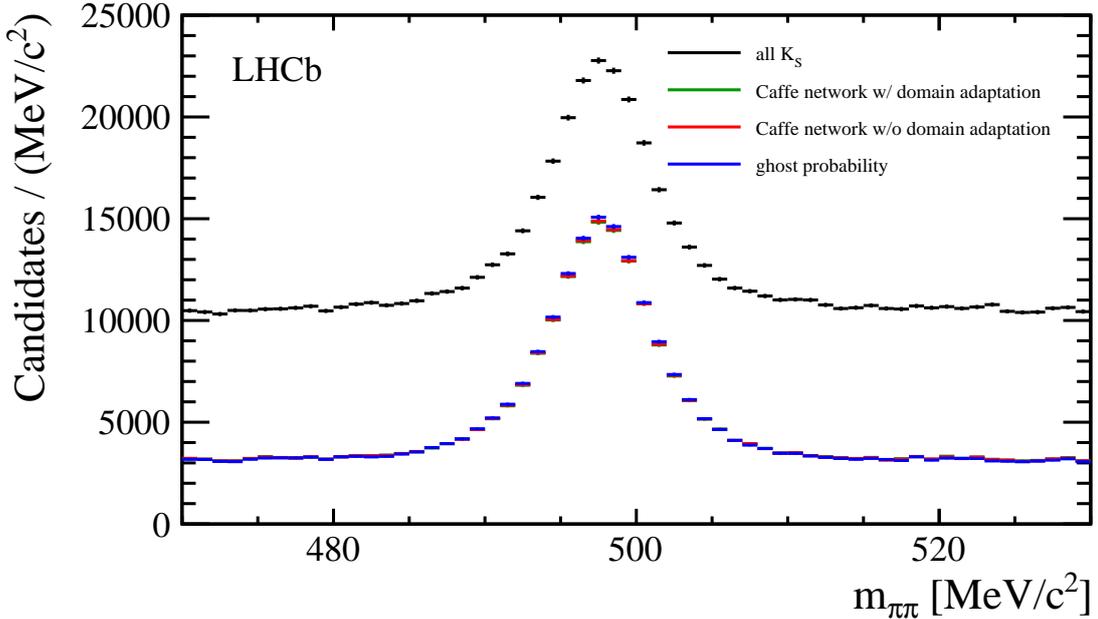


Figure 9: Comparison of the K_S^0 selection with the nominal ghost probability (blue), a network trained with the Caffe package (red), and a network trained with domain adaptation (green, covered by red). The working points of the Caffe networks are chosen to retain the same numbers of candidates. Hardly any performance difference is visible.

Table 1: Callgrind benchmark comparisons of different activation functions. Fields with n/a have not been evaluated or are not available with AVX intrinsics. The activation function used by the ghost probability is marked with (*).

function	default compiler options	AVX vectorisation by hand
\tanh	19,355,124,355	n/a
$\frac{1}{1+e^{-x}}$	21,140,125,632	n/a
$\frac{x}{\sqrt{1+x^2}}$ (*)	415,121,741	195,121,939
$\frac{x}{1+ x }$	395,121,798	195,104,759
$\max(0, x)$	155,095,875	115,095,891

281 Acknowledgements

282 We express our gratitude to our colleagues in the LHCb collaboration who provided sugges-
283 tions and helped in the implementation, especially Angelo Di Canto, Helge Voss, Manuel
284 Schiller, Gerhard Raven, Chris Jones; from the Electronic Vision(s) group (Kirchhoff-
285 Institute for Physics, Ruprecht-Karls-Universität Heidelberg, Heidelberg, Germany) Eric
286 Müller; and from the LHC Physics and New Particles group (Institute for Theoretical
287 Physics, Ruprecht-Karls-Universität Heidelberg, Heidelberg, Germany) Johann Brehmer.

288 **References**

- 289 [1] J. Weidendorfer, M. Kowarschik, and C. Trinitis, *A Tool Suite for Simulation*
 290 *Based Analysis of Memory Access Behavior*, in *Computational Science - ICCS*
 291 *2004, 4th International Conference, Kraków, Poland, June 6-9, 2004, Proceed-*
 292 *ings, Part III* (M. Bubak, G. D. van Albada, P. M. A. Sloot, and J. Dongarra,
 293 eds.), vol. 3038 of *Lecture Notes in Computer Science*, pp. 440–447, Springer, 2004.
 294 <http://www.valgrind.org/docs/pubs.html>.
- 295 [2] LHCb collaboration, R. Aaij *et al.*, *Measurement of the track reconstruction efficiency*
 296 *at LHCb*, JINST **10** (2015) P02007, [arXiv:1408.1251](https://arxiv.org/abs/1408.1251).
- 297 [3] J. Brehmer, J. Albrecht, and P. Seyfert, *Ghost probability: an efficient tool to remove*
 298 *background tracks*, <https://cds.cern.ch/record/1478372>. LHCb internal note LHCb-
 299 INT-2012-025.
- 300 [4] LHCb collaboration, R. Aaij *et al.*, *Observation of $B_c^+ \rightarrow J/\psi D_s^+$ and $B_c^+ \rightarrow J/\psi D_s^{*+}$*
 301 *decays*, Phys. Rev. **D87** (2013) 112012, [arXiv:1304.4530](https://arxiv.org/abs/1304.4530).
- 302 [5] A. Hoecker *et al.*, *TMVA: Toolkit for Multivariate Data Analysis*, PoS **ACAT** (2007)
 303 040, [arXiv:physics/0703039](https://arxiv.org/abs/physics/0703039).
- 304 [6] M. Schiller, H. Voss, and L. Moneta, *fast tanh implementation*, ROOT-7054.
- 305 [7] Wikipedia. Sigmoid function, 2015.
- 306 [8] P. Seyfert and H. Voss, *TActivation... implementations*, ROOT-7062.
- 307 [9] C. M. Bishop, *Pattern recognition and machine learning*, Information science and
 308 statistics, Springer, New York [u.a.], 10. (corrected at 8th printing) ed., 2009; J.-H.
 309 Zhong *et al.*, *A program for the Bayesian Neural Network in the ROOT framework*,
 310 Comput. Phys. Commun. **182** (2011) 2655, [arXiv:1103.2854](https://arxiv.org/abs/1103.2854).
- 311 [10] M. Heß, *Multiple scattering in track reconstruction*,
 312 <https://cds.cern.ch/record/1957764>. LHCb internal note LHCb-INT-2014-043.
- 313 [11] L. Beaucourt, E. Tournfier, M. N. Minard, and J. F. Marchand, *$B^0 \rightarrow K^* \gamma(e^+ e^-)$*
 314 *and $B_s^0 \rightarrow \phi \gamma(e^+ e^-)$ analysis status*, in *2nd Radiative decays @LHCb Workshop*
 315 (A. Oyanguren Campos *et al.*, eds.), 2015. <https://indico.cern.ch/event/375424/>.
- 316 [12] LHCb collaboration, R. Aaij *et al.*, *Measurements of the B^+ , B^0 , B_s^0 meson and Λ_b^0*
 317 *baryon lifetimes*, JHEP **04** (2014) 114, [arXiv:1402.2554](https://arxiv.org/abs/1402.2554).
- 318 [13] W. Hulsbergen, *LHCBPS-1414*, <https://its.cern.ch/jira/browse/LHCBPS-1414>.
- 319 [14] Y. Ganin and V. Lempitsky, *Unsupervised Domain Adaptation by Backpropagation*,
 320 ArXiv e-prints (2014) [arXiv:1409.7495](https://arxiv.org/abs/1409.7495).

- 321 [15] Y. Jia *et al.*, *Caffe: Convolutional Architecture for Fast Feature Embedding*,
322 [arXiv:1408.5093](https://arxiv.org/abs/1408.5093).
- 323 [16] V. Vanhoucke, A. Senior, and M. Z. Mao, *Improving the speed of neural networks on*
324 *CPUs*, <https://research.google.com/pubs/archive/37631.pdf>.
- 325 [17] P. Seyfert, *github project TMVA-MLP*, <https://github.com/pseyfert/tmva-mlp>.
- 326 [18] P. Seyfert, *CPU performance comparison of activation functions*,
327 https://twiki.cern.ch/twiki/bin/view/Main/PaulSeyfert?forceShow=1#activation_function.