

Fast neural-net based fake track rejection

Paul and others (Stephen Farry, Sascha Stahl, Michel DeCian, who else?)

Abstract

A neural-network based algorithm to identify fake tracks in the LHCb pattern recognition allows to use the same track reconstruction online (i.e. the second stage of the software trigger) and offline. A low fake rate reduces the load on the particle identification and the combinatorics of decay reconstructions. At the same time, the algorithm fits into the CPU time budget of the farms.

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. The larger timing budget, available in the trigger, and a reduced computing resource demand of the event reconstruction result in the convergence of the online and offline track reconstruction. The same performance of the track reconstruction and PID are achieved online and offline. This offers the opportunity to optimise the event selection in the trigger with stronger constraints and including the hadronic PID. It additionally increases selection efficiencies and purity and reduces systematic uncertainties.

A key ingredient to achieve this ambitious goal is the reduction of the fake track rate prior to the particle identification 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. A selection requirement at the nominal working point of 0.4 is implied¹ when the ghost probability is referred to as a selection requirement.

¹this corresponds to a fake track retention of 40%

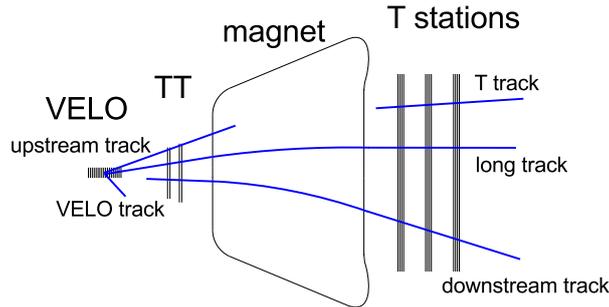


Figure 1: Illustration of the tracking system of LHCb, starting from the VELO around the collision point on the left, particles pass the TT, are deflected in the magnetic field of the dipole magnet and then detected in the T station detector (IT and OT). Different track types are reconstructed by different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

38 1.1 Track reconstruction

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

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

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

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

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

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

74 1.2 Previous works

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

80 The 22 input variables to the old ghost probability are the track fit χ^2 , and the individual
81 contributions $\chi_D^2, \chi_U^2, \chi_M^2$ and the corresponding degrees of freedom; the numbers of hits
82 on the track in each tracking detector; the reconstructed track p_T and pseudorapidity;
83 the track is interpolated through the detector to count in how many active strips/straws
84 no hit is observed although the track passes through it (“expected hits”); and finally the
85 occupancies of all tracking detectors.

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

88 1.3 Training sample

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

- 93 • 25 ns bunch spacing, $\nu = 1.6$
- 94 • 25 ns bunch spacing, $\nu = 1.9$
- 95 • 50 ns bunch spacing, $\nu = 1.6$
- 96 • 50 ns bunch spacing, $\nu = 2.7$

97 The scenarios differ, for what concerns the track reconstruction, significantly in detector
98 occupancy and spillover in the Outer Tracker. That may lead to different behaviours
99 of fake track reconstruction and require different network trainings for different running
100 conditions. The necessity for having different network trainings is assessed by training

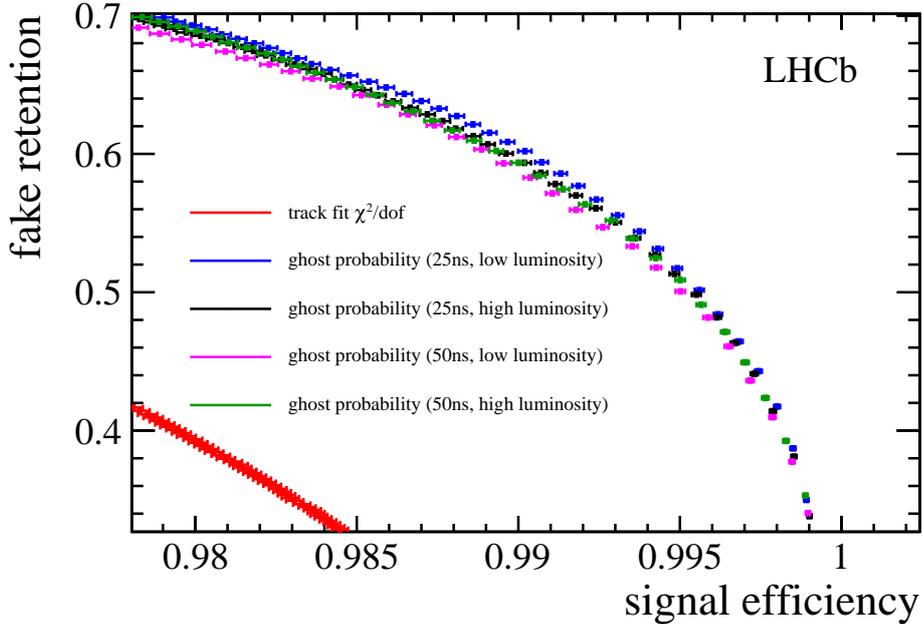


Figure 2: ROC curves for fake track discriminating neural networks, trained for different LHC running conditions, evaluated for 25 ns bunch spacing, $\nu = 1.6$. The error bars are statistical uncertainties only.

101 networks for each of the running conditions, with all other training parameters fixed,
 102 and evaluating the networks on one of the samples and their discriminating powers are
 103 compared. Figure 2 shows the ROC curves for the 25 ns sample at low pile-up. The
 104 discrimination powers of the four networks do not largely differ and thus for simplicity only
 105 a single network (trained at the favoured scenario of 25 ns bunch spacing at low pile-up) is
 106 deployed.

107 1.4 Network architecture tuning

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

- 109 • equipped with a root file interface for the training, which is the common data file
 110 format in LHCb software,
- 111 • commonly known in LHCb (ensuring future maintainability),
- 112 • and provides code generation for the trained network such that the network can be
 113 integrated into any C++ code without creating dependency on external libraries.

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

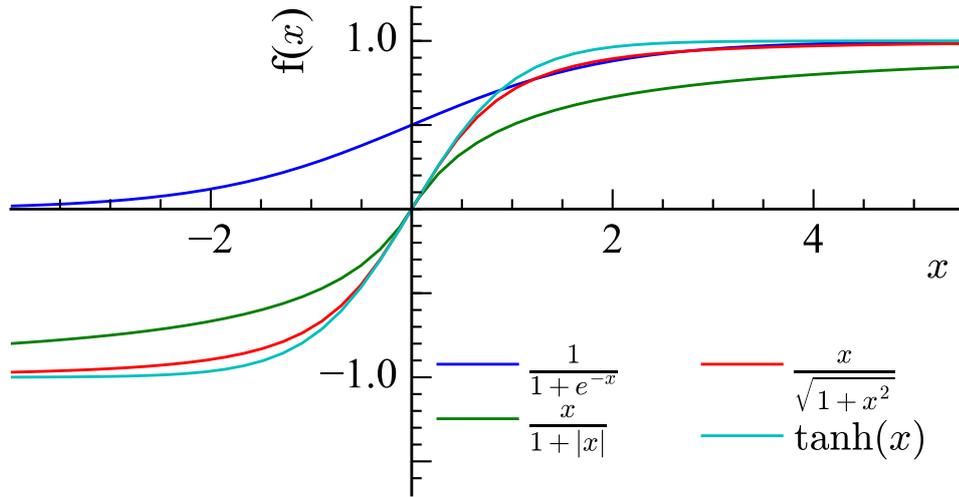


Figure 3: Functional shape of sigmoid functions.

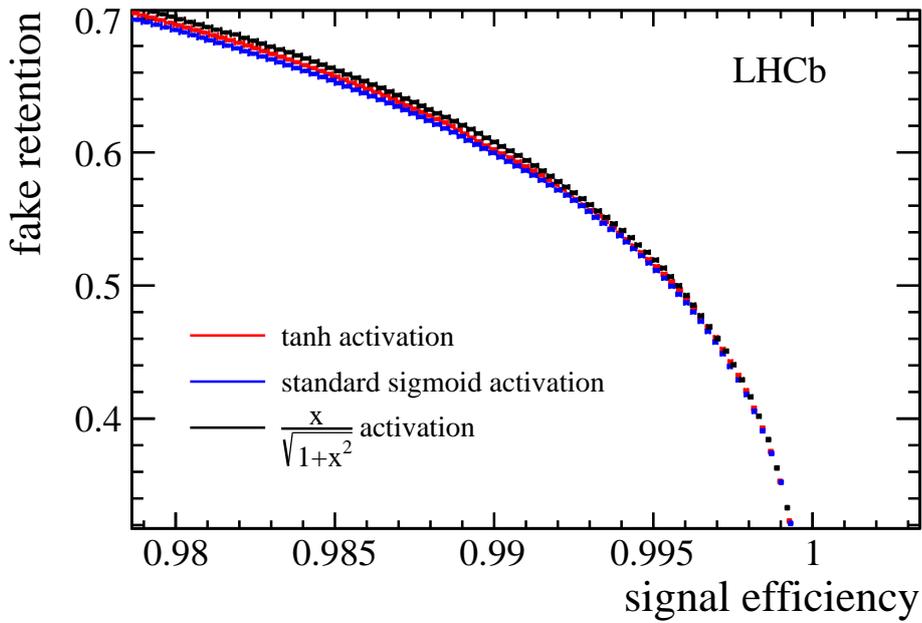


Figure 4: ROC curves for fake track discriminating neural networks, using different activation functions. A very small physics performance improvement is observed when changing from the TMVA standard functions to $\frac{x}{\sqrt{1+x^2}}$.

118 Of the tested functions $\frac{x}{\sqrt{1+x^2}}$ is the fastest to compute, while no significant physics
 119 performance difference is expected given the similar functional shape, see Fig. 3. No
 120 performance difference is observed in Fig. 4. Therefore, it is chosen as activation function.

121 The ghost probability is a classification problem, and thus cross entropy [9] is chosen
122 as loss function in the network training. With respect to the run I implementation of the
123 ghost probability, this contributes to the physics performance improvement. Technically
124 the output layer activation function when training with cross entropy loss must be $\frac{1}{1-e^{-x}}$
125 in TMVA, it is removed from the network after the training and replaced by a custom
126 calibration as described in Sect. 1.6).

127 1.5 Variable selection

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

133 1.6 Output transformation

134 To ease the usage of the ghost probability, a transformation of the network is applied. A
135 probability integral transform² is obtained as a linear spline fit to the cumulative network
136 response for fake tracks in simulated events. The discriminating behaviour of any classifier
137 is invariant under monotonous transformations and so is the physics performance under
138 the probability integral transform. Motivations for this transformation are primarily to
139 give a physical interpretation to the response: rejecting tracks with a ghost probability of
140 larger than $x\%$ will retain $x\%$ of all fake tracks.

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

144 1.7 Category classifiers

145 Fake tracks produced by different pattern recognition algorithms might have different
146 track properties. It might therefore be beneficial to train separate neural networks for
147 the two main track reconstruction algorithms at LHCb. On simulated events, the physics
148 performance of two separate networks does not differ from the physics performance of a
149 single network. Similarly, different networks for different T station regions have been tested
150 (one for tracks in the OT, IT, and the overlap region), without significant performance
151 gain.

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

²also referred to as “flattening” or “rarity transformation”

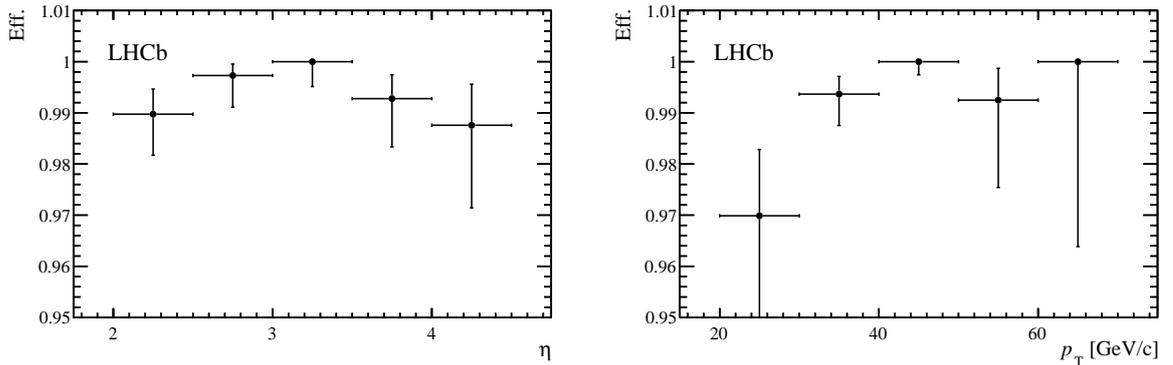


Figure 5: Efficiency for $Z \rightarrow \mu\mu$ tracks to pass the ghost probability, in data from 2015 at a bunch spacing of 50 ns.

2 Validation

The data taking strategy of LHCb in Run II involves the application of the same track reconstruction in the software trigger as in the offline data processing. This goal can only be achieved by applying the ghost probability in the software trigger to reduce the rate of fake tracks entering the particle identification and combinatorics of decay reconstructions.

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

The latest update to the track fit χ^2 computation, as used since 25 ns data taking in 2015, was applied in all validations by refitting all tracks prior to computing the ghost probability.

2.1 High momentum tracks (data from 2015, 50 ns bunch spacing)

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

2.2 Electron reconstruction (data from run I, 50 ns bunch spacing)

Electrons are more challenging to reconstruct than the standard candles. At the same time, it can be expected that the response of the ghost probability for electrons differs

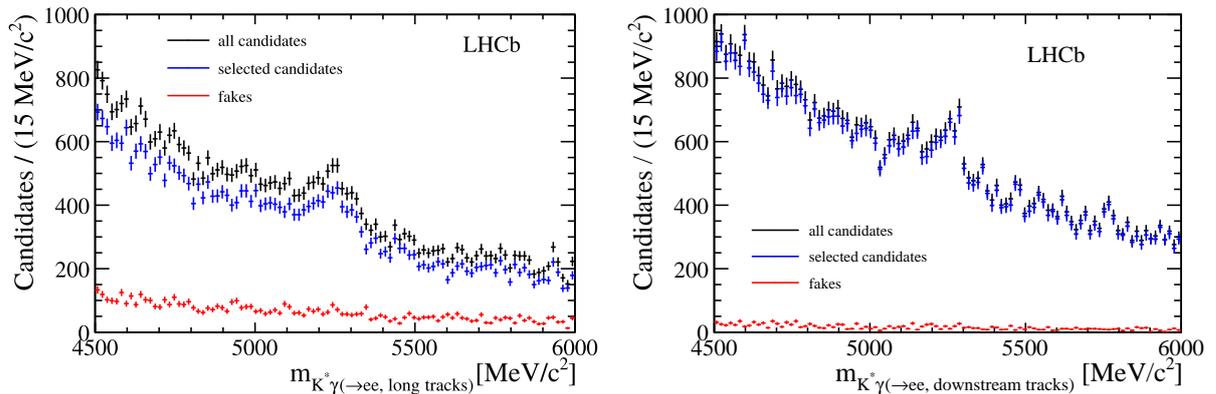


Figure 6: $B_s^0 \rightarrow K^* \gamma$ candidates, where the photon is reconstructed as pair of electron tracks, using long tracks (left) and using downstream tracks (right). The candidates prior to a cut on the ghost probability of the electron tracks are shown in black, those passing the ghost probability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

177 from that for other particles as electrons undergo more multiple scattering.

178 It is assumed that the reconstruction of converted photons as electron pair is the most
 179 vulnerable channel since the photon conversion can happen “late” in the VELO leaving
 180 only few hits, the e^+e^- pair has a small opening angle which could lead to hit ambiguities
 181 in the VELO pattern reconstruction, and analyses of channels $B_s^0 \rightarrow K^* \gamma$ are anyhow
 182 so-called rare decays which immediately suffer from efficiency loss.

183 The “early data” of 2015 does not correspond to enough integrated luminosity to obtain
 184 a satisfying estimation of the consequences of a cut on the ghost probability. The tracks of
 185 $B_s^0 \rightarrow K^* \gamma$ candidates from Run I³ are refitted using the track fit configuration as used in
 186 2015 and the ghost probability is evaluated. The invariant mass spectrum shown in Fig. 6
 187 shows candidates without the application of a cut on the ghost probability, those passing,
 188 and those failing; both for converted photons reconstructed as pair of downstream tracks
 189 and as pair of long tracks. To the statistical precision of this test, no signal loss is visible.

190 2.3 Validation with 25 ns data from 2015

191 A cut on the ghost probability is included in the standard track reconstruction since data
 192 taking at 25 ns bunch spacing started. To investigate the behaviour of the ghost probability
 193 in real data with 25 ns bunch spacing, events are re-reconstructed without a cut on the
 194 ghost probability. Under the assumption, that most K_s^0 are part of the underlying event
 195 and most triggered events containing K_s^0 would have been triggered without those K_s^0 , K_s^0
 196 are used as probe of the ghost probability.

197 The invariant mass spectrum of K_s^0 candidates after re-reconstruction is shown in Fig. 7
 198 for both long tracks and downstream tracks. In either case, the background contribution

³using a simplified version of the selection presented in [10] – without Bremsstrahlung correction

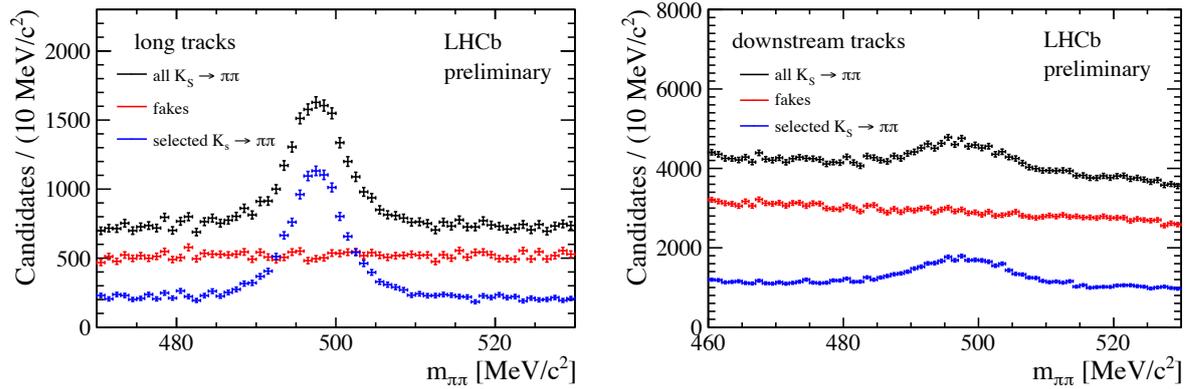


Figure 7: StdLooseKS invariant mass spectrum for events reconstructed without using the ghost probability in the track selection; using long tracks (left) and downstream tracks (right). Candidates for which at least one track fails the default ghost probability requirement are shown in red and do not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

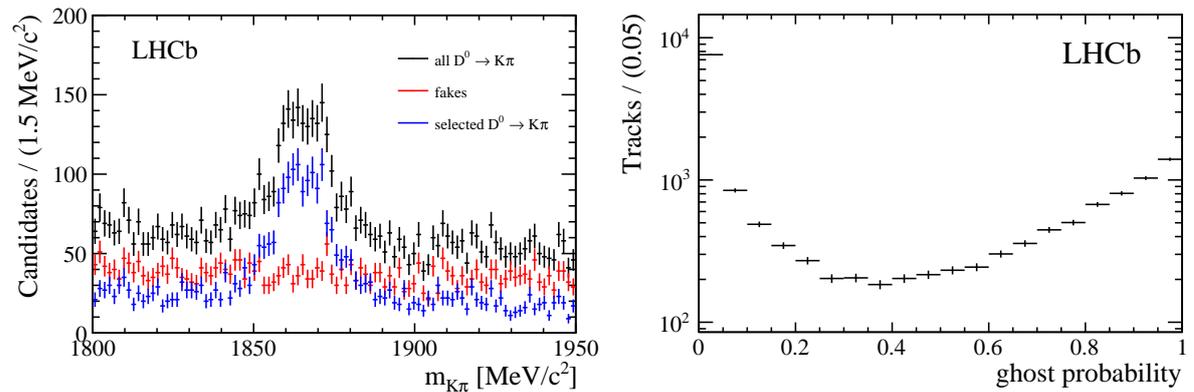


Figure 8: StdLooseD02KPi invariant mass spectrum for events reconstructed without using the ghost probability in the track selection (left). Candidates for which at least one track fails the default ghost probability requirement are shown in red and do not exhibit a signal contamination. The remaining candidates are shown in blue. The ghost probability distribution for the kaon and pion tracks is shown on the right.

199 is largely reduced when rejecting events where at least one of the tracks has a ghost
 200 probability of larger than 0.4. There is no signal visible in the events rejected. It is
 201 concluded that no physics signal is lost due to the application of the ghost probability.⁴

202 The same test with $D \rightarrow K\pi$ decays is shown in Fig. 8. To ensure that the sample
 203 is not biased towards candidates passing the ghost probability due to the online event
 204 selection, the ghost probability spectrum is shown in Fig. 8 b), where no step from such a
 205 selection is visible at 0.4.

⁴within the statistical sensitivity of the test; in the kinematic spectrum of the selected K_S^0

2.4 Decay time acceptance (simulated events)

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

To evaluate a possible decay time bias of the ghost probability, for each reconstructed particle in simulated events with a $b\bar{b}$ production, the average of the true decay time of their ancestor particles is determined. When rejecting tracks which fail the ghost probability criterion, the average ancestor decay time changes by $(1.5 \pm 2.0) \times 10^{-15}$ s. This is smaller than the statistical sensitivity of this test, and smaller than the systematic uncertainty to which the lifetime bias of the LHCb reconstruction is known [11].

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 [12].

The current training is purely based on simulated events, the domain adaptation approach from [13] 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 [14], which is used in [13]. 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 [15] lead to tenfold improvement of the neural network implementation [16] once using AVX intrinsic commands. This approach has not been followed up to ensure platform independence of the ghost probability.

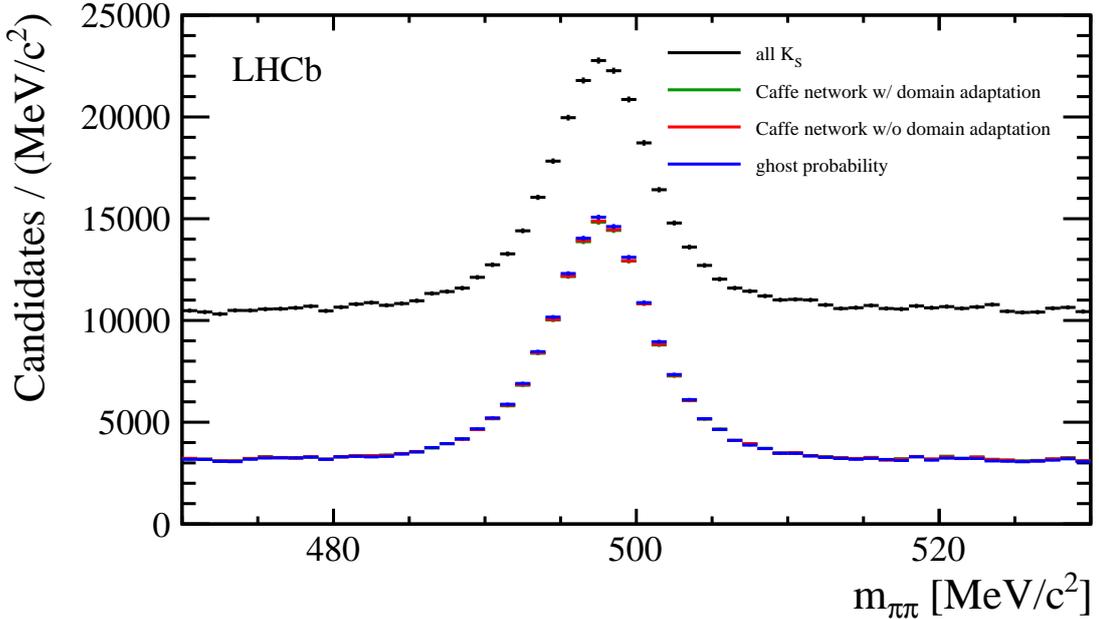


Figure 9: Comparison of the K_s^0 selection with the nominal ghost probability (blue), a network trained with the Caffe package (red), and a network trained with domain adaptation (green, covered by red). The working points of the Caffe networks are chosen to retain the same numbers of candidates. Hardly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

246 The current activation function in the neural network is $\frac{x}{\sqrt{1+x^2}}$. The rectified linear
 247 unit $\max(0, x)$ or $\frac{x}{1+|x|}$ are expected to be even faster, as listed in Tab. 1 (from [17]).

248 4 Conclusion

249 The ghost probability is a super awesome tool which allows LHCb to do super good physics,
 250 and we know we can use it for everything and it makes our software trigger fast, which
 251 means we can do more fancy stuff with the CPU resources.

252 References

- 253 [1] J. Weidendorfer, M. Kowarschik, and C. Trinitis, *A Tool Suite for Simulation*
254 *Based Analysis of Memory Access Behavior*, in *Computational Science - ICCS*
255 *2004, 4th International Conference, Kraków, Poland, June 6-9, 2004, Proceed-*
256 *ings, Part III* (M. Bubak, G. D. van Albada, P. M. A. Sloot, and J. Dongarra,
257 eds.), vol. 3038 of *Lecture Notes in Computer Science*, pp. 440–447, Springer, 2004.
258 <http://www.valgrind.org/docs/pubs.html>.
- 259 [2] LHCb collaboration, R. Aaij *et al.*, *Measurement of the track reconstruction efficiency*
260 *at LHCb*, JINST **10** (2015) P02007, [arXiv:1408.1251](https://arxiv.org/abs/1408.1251).
- 261 [3] J. Brehmer, J. Albrecht, and P. Seyfert, *Ghost probability: an efficient tool to remove*
262 *background tracks*, <https://cds.cern.ch/record/1478372>. LHCb internal note LHCb-
263 INT-2012-025.
- 264 [4] LHCb collaboration, R. Aaij *et al.*, *Observation of $B_c^+ \rightarrow J/\psi D_s^+$ and $B_c^+ \rightarrow J/\psi D_s^{*+}$*
265 *decays*, Phys. Rev. **D87** (2013) 112012, [arXiv:1304.4530](https://arxiv.org/abs/1304.4530).
- 266 [5] A. Hoecker *et al.*, *TMVA: Toolkit for Multivariate Data Analysis*, PoS **ACAT** (2007)
267 040, [arXiv:physics/0703039](https://arxiv.org/abs/physics/0703039).
- 268 [6] M. Schiller, H. Voss, and L. Moneta, *fast tanh implementation*, ROOT-7054.
- 269 [7] Wikipedia. Sigmoid function, 2015.
- 270 [8] P. Seyfert and H. Voss, *TActivation... implementations*, ROOT-7062.
- 271 [9] C. M. Bishop, *Pattern recognition and machine learning*, Information science and
272 statistics, Springer, New York [u.a.], 10. (corrected at 8th printing) ed., 2009; J.-H.
273 Zhong *et al.*, *A program for the Bayesian Neural Network in the ROOT framework*,
274 Comput. Phys. Commun. **182** (2011) 2655, [arXiv:1103.2854](https://arxiv.org/abs/1103.2854).
- 275 [10] L. Beaucourt, E. Tournfier, M. N. Minard, and J. F. Marchand, *$B^0 \rightarrow K^* \gamma(e^+e^-)$*
276 *and $B_s^0 \rightarrow \phi \gamma(e^+e^-)$ analysis status*, in *2nd Radiative decays @LHCb Workshop*
277 (A. Oyanguren Campos *et al.*, eds.), 2015. <https://indico.cern.ch/event/375424/>.
- 278 [11] LHCb collaboration, R. Aaij *et al.*, *Measurements of the B^+ , B^0 , B_s^0 meson and Λ_b^0*
279 *baryon lifetimes*, JHEP **04** (2014) 114, [arXiv:1402.2554](https://arxiv.org/abs/1402.2554).
- 280 [12] W. Hulsbergen, *LHCbps-1414*, <https://its.cern.ch/jira/browse/LHCbps-1414>.
- 281 [13] Y. Ganin and V. Lempitsky, *Unsupervised Domain Adaptation by Backpropagation*,
282 ArXiv e-prints (2014) [arXiv:1409.7495](https://arxiv.org/abs/1409.7495).
- 283 [14] Y. Jia *et al.*, *Caffe: Convolutional Architecture for Fast Feature Embedding*,
284 [arXiv:1408.5093](https://arxiv.org/abs/1408.5093).

- 285 [15] V. Vanhoucke, A. Senior, and M. Z. Mao, *Improving the speed of neural networks on*
286 *CPUs*, <https://research.google.com/pubs/archive/37631.pdf>.
- 287 [16] P. Seyfert, *github project TMVA-MLP*, <https://github.com/pseyfert/tmva-mlp>.
- 288 [17] P. Seyfert, *CPU performance comparison of activation functions*,
289 https://twiki.cern.ch/twiki/bin/view/Main/PaulSeyfert?forceShow=1#activation_function.