



sPH-HF-2023-001

v0.1

January 13, 2023

# $D^0 \rightarrow K^- \pi^+$ modelling and selection in min-bias Au+Au simulations at sPHENIX

Cameron Dean

*Massachusetts Institute of Technology*

## **Abstract**

The reconstruction and selection of  $D^0 \rightarrow K^- \pi^+$  in simulated Au+Au collisions with the sPHENIX detector is presented. Approximately 22 million minimum bias events were generated at  $\sqrt{s_{NN}} = 200$  GeV using the HIJING event generator, corresponding to roughly 20 minutes of data taking at a collision rate of 15 kHz.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Simulation</b>	<b>2</b>
<b>3</b>	<b>Toolkit</b>	<b>3</b>
3.1	KFParticle . . . . .	3
3.2	DecayFinder . . . . .	4
3.3	HFTrackEfficiency . . . . .	5
<b>4</b>	<b>Selection</b>	<b>5</b>
<b>5</b>	<b>Models</b>	<b>9</b>
5.1	Signal . . . . .	9
5.2	Background . . . . .	10
5.2.1	Additional background sources . . . . .	11
<b>6</b>	<b>Results</b>	<b>11</b>
<b>7</b>	<b>Conclusion</b>	<b>16</b>
	<b>Appendices</b>	<b>19</b>
<b>A</b>	<b>Alternative Fit Models</b>	<b>19</b>
A.1	Bifurcated Gaussian . . . . .	19
A.2	Double Crystal Ball . . . . .	20
A.3	Kernel density estimated backgrounds . . . . .	22
<b>B</b>	<b>Machine Learning Selection</b>	<b>24</b>
<b>C</b>	<b>sWeighting</b>	<b>30</b>

# 1 Introduction

1 The reconstruction of open-charm hadrons is a key requirement of the sPHENIX experiment  
 2 which will commence data-taking in Spring 2023. As with any new detector, the initial data-  
 3 taking period will be devoted to commissioning the experiment. This involves ensuring the  
 4 correct operation of all subsystems, the tracking, calorimeter calibration and calculation of  
 5 physics objects. Preparations for the latter requirement is the focus of this note. A detailed  
 6 commissioning timeline has been agreed by the collaboration which is publicised within the  
 7 Beam Use Proposal [1]. The timeline is given in Table 1 for completeness.

Weeks	Details
2.0	low rate, 6-28 bunches
2.0	low rate, 111 bunches, MBD L1 timing
1.0	low rate, crossing angle checks
1.0	low rate, calorimeter timing
4.0	medium rate, TPC timing, optimization
2.0	full rate, system test, DAQ throughput
<b>12.0</b>	<b>Total</b>

Table 1: Timeline for sPHENIX commissioning period in 2023, the first year of operation. Copied from the sPHENIX 2022 beam use proposal [1].

8 A challenging aspect of the commissioning period will be the reconstruction of  $D^0 \rightarrow$   
 9  $K^- \pi^+$  as it is both reasonably rare and the decay products often have a transverse mo-  
 10 mentum ( $p_T$ ) less than 1 GeV. This channel is also key to realising the  $b$ -physics program  
 11 of sPHENIX as the separation of prompt and non-prompt open-charm decays is used to  
 12 tag probable  $b$ -hadrons [2, 3, 4]. The number of  $D^0 \rightarrow K^- \pi^+$ ,  $N_{D^0 \rightarrow K^- \pi^+}$ , produced is  
 13 proportional to the integrated luminosity,  $\mathcal{L}$ , and is roughly given by

$$N_{D^0 \rightarrow K^- \pi^+} \approx \mathcal{L} \sigma_{c\bar{c}} 2N_{\text{coll}} f_{D^0} \text{BF}(D^0 \rightarrow K^- \pi^+) \varepsilon_{\text{acc}}^2 \varepsilon_{\text{track}}^2 \varepsilon_{\text{sel}} \quad (1)$$

14 where  $\sigma_{c\bar{c}}$  is the production cross section for charm-quarks (about  $1/60^{\text{th}}$  of the total  
 15 inelastic cross-section [5, 6]),  $N_{\text{coll}}$  is the number of binary collisions,  $f_{D^0}$  is the  $D^0$  fragmen-  
 16 tation fraction (about 40% [7]),  $\text{BF}(D^0 \rightarrow K^- \pi^+)$  is the  $D^0 \rightarrow K^- \pi^+$  branching fraction  
 17 (about 4% [8]),  $\varepsilon_{\text{acc}}$  is the geometrical acceptance efficiency (about 23%, see Section 3.2),  
 18  $\varepsilon_{\text{track}}$  is the tracking efficiency (about 80%, see Section 3.3) and  $\varepsilon_{\text{sel}}$  is the selection effi-  
 19 ciency. The factor of two comes from producing 2 charm quarks which are equally likely to  
 20 go through a  $D^0 \rightarrow K^- \pi^+$  decay and the squares on the acceptance and tracking efficiencies  
 21 comes about from both efficiencies applying equally to each daughter track. This means the  
 22 the expected yield of  $D^0 \rightarrow K^- \pi^+$  is given by

$$N_{D^0 \rightarrow K^- \pi^+} \approx \mathcal{L} N_{\text{coll}} \varepsilon_{\text{sel}} \times 10^{-5} \quad (2)$$

---

<sup>1</sup>Charge conjugation is implied throughout this note unless otherwise stated.

23 The aim of this note is to understand the topological and kinematic variable distributions  
 24 of  $D^0$  decays in Au+Au collisions at sPHENIX and hence optimise  $\varepsilon_{\text{sel}}$ . This note is intended  
 25 to act as a baseline reference for the initial selection of  $D^0 \rightarrow K^-\pi^+$  in the commissioning  
 26 period. Several aspects of the experiment data pipeline will be improved during the commis-  
 27 sioning period as the collaboration sees and understands the initial data, most importantly  
 28 the tracking. It is expected that the momentum resolution will undergo improvements in  
 29 the initial period and this resolution has a direct impact on the  $D^0$  resolution. Due to the  
 30 large data volume and need for a fast turnaround on physics object construction, we will  
 31 pre-select tracks that share signatures with a heavy flavor decay (as examples, tracks with  
 32 a large  $p_T$  or large distance-of-closest approach to the primary vertex) and then save these  
 33 track seeds and clusters to a smaller data file to re-run the track fitting with improved align-  
 34 ment parameters. It is important then to understand what these variables look like and  
 35 select tracks with loose enough cuts to catch tracks that may not have the best resolution  
 36 yet but the cuts need to be tight enough to reject a large portion of the background. The  
 37 reconstruction of  $D^0$  then serves as an important indicator that the sPHENIX collaboration  
 38 is ready to produce physics results.

39 The note after this introduction is arranged as follows; the general simulation setup of the  
 40 generators and detector, along with the generated statistics is detailed in Section 2, the tools  
 41 developed within sPHENIX to understand the decay topology are described in Section 3,  
 42 the base selections applied to the simulation data set is detailed in Section 4, the models  
 43 used to describe the invariant mass distributions are described in Section 5 and the results  
 44 are presented in Section 6.

45 To avoid biases, the signal was modelled using  $D^0 \rightarrow K^-\pi^+$  decays simulated in  $p+p$   
 46 collisions using PYTHIA8 while the background was modelled using HIJING events outside  
 47 of the invariant mass search window. Comparisons of the signal-to-background ratio are  
 48 made between candidates selected using direct cuts and machine-learning methods.

## 49 2 Simulation

50 This study used two different data sets. One data set consisted of  $p+p$  collisions at  $\sqrt{s} = 200$  GeV  
 51 using PYTHIA8 [9] as an event generator. PYTHIA8 was tuned to approximate the minimum  
 52 bias collision environment at RHIC [10] but, after the initial generation, each event was re-  
 53 quired to have a  $c\bar{c}$  pair produced to enrich the  $D^0 \rightarrow K^-\pi^+$  statistics. The RHIC collision  
 54 rate for  $p+p$  collisions is almost 10 MHz and some of the sPHENIX detectors integrate their  
 55 hits over a longer period than this and hence, there is an out-of-time pileup effect. This effect  
 56 was implemented as part of our simulation. A total of 50 million events were generated in  
 57 this fashion and they were used to model the  $D^0 \rightarrow K^-\pi^+$  signal shape and variable distri-  
 58 butions. This was intended to avoid biases in the final selection and fitting by not using the  
 59 same  $D^0 \rightarrow K^-\pi^+$  candidates for both selection and testing.

60 The second data set consisted of minimum bias Au+Au collisions at  $\sqrt{s_{NN}} = 200$  GeV  
 61 using HIJING [11]. Unlike the  $p+p$  sample, there were no further requirements on the  
 62 generator to enrich the  $D^0 \rightarrow K^-\pi^+$  purity. This sample was used for both background  
 63 modelling and for the final test of the selection. The background sample was taken from  
 64 a region outside the invariant mass range to ensure there is no bias in the final selection.

65 Almost 22 million events were generated which corresponds to a little over 20 minutes of  
66 data taking at the peak RHIC Au+Au collision rate of 15 kHz.

67 The detector was simulated using the GEANT4 package [12]. The simulation consists  
68 of the beam pipe and all detectors, except the event plane detector (sEPD) which is not  
69 used for this study. In the interaction region, the beam pipe consists of a beryllium section  
70 with an inner radius of 2 cm and a thickness of 762  $\mu\text{m}$ . The volume inside the beam pipe  
71 is simulated as a vacuum. The monolithic active pixel sensor vertex detector (MVTX) is  
72 simulated after the beam pipe and extends from 2.4 cm to 5.3 cm. It consists of a three  
73 layer vertexing detector with a timing resolution of approximately 5  $\mu\text{s}$  and hence sees the  
74 out of time pile-up effect. This effect is accounted for by duplicating the hits with a time  
75 delay in GEANT4. The next detector radially is the intermediate tracker (INTT) which  
76 consists of a two-layer silicon strip detector with a timing resolution capable of resolving the  
77 10 MHz collision rate. The INTT sits between 7.2 cm and 10.3 cm, radially. A compact time  
78 projection chamber (TPC) sits beyond the INTT up to a radius of 80 cm and provides the  
79 momentum measurement for tracks in sPHENIX. The TPC integrates over a period of 34  $\mu\text{s}$   
80 and also sees the out-of-time pileup effect. The final tracking detector is the TPC outer  
81 tracker (TPOT) which consists of 8 GEM detectors below the TPC in the vertical direction.  
82 This acts as a final measurement point in the tracking to account for the distortions and drift  
83 in the TPC tracks. As well as the tracking detectors, there are two hadronic calorimeters  
84 and an electromagnetic calorimeter in the simulation. As they are not used to identify the  
85  $D^0$  candidates, they are not described in this note. Between the two hadronic calorimeters,  
86 there is a 1.5 T superconducting solenoidal magnet which is also simulated. The tracking is  
87 performed using A Common Tracking Software (ACTS) [13, 14].

## 88 3 Toolkit

### 89 3.1 KFParticle

90 The candidate reconstruction was performed by KFPARTICLE, originally designed by the  
91 CBM collaboration [15] and adapted to be used within the Fun4All framework [16]. The  
92 internal logic of the package has been largely unchanged since its previous description in the  
93 sPHENIX framework [17] except for small bug fixes or additional output variables for users.  
94 The main update since the previous note is the addition of a “decay descriptor” to simplify  
95 the user interface.

96 The decay descriptor is a string that users write to specify the decay topology and is  
97 parsed by the top interface class before any event processing. If the descriptor can not be  
98 understood by the parser, a warning will be raised and the module will not be added to the  
99 node tree although Fun4All will still run. The decay parser checks each particle the user  
100 specifies against the particle database found in ROOT’s TDatabasePDG to ensure it exists.  
101 If the particle does not exist, this will be written to the user. Similarly, if the parser cannot  
102 interpret the charge of a track, the user will also be notified of the offending track. All other  
103 instances of the parser not understanding the string will be written as a standard warning.  
104 An example decay descriptor is as follows

```
105 [B+ -> {D0bar -> K^- pi^+} pi^+]cc
```

106 This is interpreted as a  $B^+$  hadron decaying to an intermediate  $D^0$  and an isolated  
 107 charged pion. The mother is always written to the left of a  $\rightarrow$  and intermediate decays  
 108 are always contained within  $\{ \}$  braces. The charge of a track is written as either +, - or  
 109 0 directly after a caret ( $\wedge$ ). If you wish to also search for the charge conjugate decay, the  
 110 entire descriptor must be contained with square braces and appended with `cc` or `CC` for  
 111 charge-conjugate. From this, `KFPARTICLE` knows how many intermediate decays there are,  
 112 how many tracks are in each intermediate state and final state as well as the charges of all  
 113 tracks. It also knows if it needs to bring the intermediate states back to a common vertex  
 114 or associate them to any other final state tracks.

## 115 3.2 DecayFinder

116 `DECAYFINDER` is a new package that runs over the truth record of an event. It searches for  
 117 decays specified by the user and accepts input ranges for the final state particles pseudorapid-  
 118 ity,  $\eta$ , and  $p_T$ . It effectively measures the geometrical acceptance of a decay. `DECAYFINDER`  
 119 also uses a decay descriptor with the same logic as `KFPARTICLE`, so users can just repeat the  
 120 string they wrote previously. `DECAYFINDER` begins by looking for the `HEPMC2` record [18]  
 121 on the node tree, if this doesn't exist then it will fall back to using the `GEANT4` truth record.  
 122 The tool loops over the truth container to find the required mother and, when this is found,  
 123 the decay products will be analysed. There is an internal list of resonances that will be  
 124 further analysed if seen in the mothers decay chain. For example, a  $\phi(1020)$  resonance often  
 125 decays to two kaons. If a user is looking for two kaons in the final state, then the  $\phi(1020)$  will  
 126 be studied further. This resonance would not automatically be seen by `sPHENIX` without  
 127 reconstructing the dikaon pair first. If the  $\phi(1020)$  is specified in the decay descriptor, it  
 128 will be removed from the internal resonance list and treated as an integral part of the decay  
 129 chain.

130 It is possible to limit the decay volume of some event generators such as `PYTHIA8`. If  
 131 this occurs, the `HEPMC2` record will not contain the full information of the event. In this  
 132 case, `DECAYFINDER` can detect that the decay volume was limited and switch to searching  
 133 the `GEANT4` truth record when this boundary is discovered<sup>2</sup>. When `DECAYFINDER` detects  
 134 a final state track ( $e^\pm$ ,  $\mu^\pm$ ,  $\pi^\pm$ ,  $K^\pm$ ,  $p$  or  $\bar{p}$ ), it calculates the track's  $\eta$  and  $p_T$  and compares  
 135 this to the user specified values<sup>3</sup>. If all tracks pass the requirement, `DECAYFINDER` tags  
 136 the decay as reconstructable. The barcodes, embedding IDs and PDG IDs of all particles  
 137 in a reconstructable decay can be written to the node tree for further analysis. When  
 138 `DECAYFINDER` finishes, it can write a report of how many decays were generated, how  
 139 many had at least one track fail the  $p_T$  requirement or had at least one track fail the  $\eta$   
 140 requirement or had at least one track fail both the  $p_T$  and  $\eta$  requirements and finally how  
 141 many decays fell within the required acceptance. `DECAYFINDER` is also capable of triggering  
 142 on events that have all particles within the `sPHENIX` acceptance.

<sup>2</sup>This boundary appears as a null pointer when calling a particle's end vertex in `HEPMC2`

<sup>3</sup>The default values are  $p_T \geq 0.2 \text{ GeV}$  and  $|\eta| \leq 1.1$

### 143 3.3 HFTrackEfficiency

144 HFTRACKEFFICIENCY uses the output of DECAYFINDER to match final state tracks to the  
145 decay products that DECAYFINDER claims are all in the required acceptance region. When  
146 each trackable particle is found in the truth record, the reconstructed track map is iterated  
147 over and the 3-momentum of the truth and reconstructed objects is compared. If all three  
148 values match within a specified limit<sup>4</sup> then HFTRACKEFFICIENCY counts that particle as  
149 reconstructed. The user can specify that they would like to have the results of the search  
150 output to an nTuple. This nTuple contains the truth momenta, PID and  $\eta$  for each trackable  
151 particle as well as the reconstructed momenta of the track, if it exists. There is also a boolean  
152 flag to say whether that track was or was not reconstructed. The mothers true momenta,  
153 PID and  $\eta$  are also written to the file along with the reconstructed mass as seen by ACTS,  
154 assuming all tracks were reconstructed. This files gives a direct look at the tracking efficiency  
155 for heavy flavor particles.

156 As with DECAYFINDER, HFTRACKEFFICIENCY is also capable of triggering when all  
157 final state particles are reconstructed. As well as this triggering, it can also write the recon-  
158 structed tracks back to the node tree in a subset of the track map. This is only done when  
159 all tracks are reconstructed for a decay and so allows users to build a decay without any  
160 selections which is useful for studying the decays kinematic distributions.

## 161 4 Selection

162 The  $D^0$  candidates were selected using kinematic and topological variable cuts. The cuts  
163 are applied to the daughter tracks, secondary vertex and reconstructed mother candidates  
164 and requires knowledge of the primary vertex (PV). The tracks are selected based on their  
165 transverse momentum ( $p_T$ ), track  $\chi^2$  per number of degrees of freedom and minimum distance  
166 of closest approach with respect to each reconstructed primary vertex (IP). The secondary  
167 vertex (SV) is selected by making pairs of tracks and measuring their distance of closest  
168 approach with respect to each other (DCA) and the  $\chi^2$  per number of degrees of freedom  
169 of the reconstructed SV. The mother candidates are selected by requiring they lie within an  
170 invariant mass range ( $m_{K-\pi^+}$ ), they have a minimum IP  $\chi^2$  with respect to its selected PV,  
171 a minimum cosine of the angle between the flight direction and mother momentum vector  
172 (DIRA), and  $p_T$ . Further variables were also studied such as the quality of the separation of  
173 the primary and secondary vertices (the flight distance  $\chi^2$ ) and the quality of the track IP.  
174 These variables were found to have little separation power and improving their calculations  
175 will be a task to undertake before the initial sPHENIX data taking. A visual description of  
176 the variables is given in Figure 1.

177 To ensure the final selection was unbiased, signal and background samples were taken from  
178 events that were not used in the final simulation. The signal events were taken from a simu-  
179 lated sample of minimum bias  $p+p$  collisions using PYTHIA8 which were filtered at the gener-  
180 ator level to ensure each event contained a  $c\bar{c}$  pair. DECAYFINDER was run on each event  
181 to select  $D^0 \rightarrow K^-\pi^+$  decays and required the minimum track  $p_T$  was greater than 0.16 GeV  
182 with a track pseudorapidity ( $\eta$ ) between  $-1.6 \leq \eta \leq 1.6$ . DECAYFINDER was required to

---

<sup>4</sup>The default matching limit is 5%

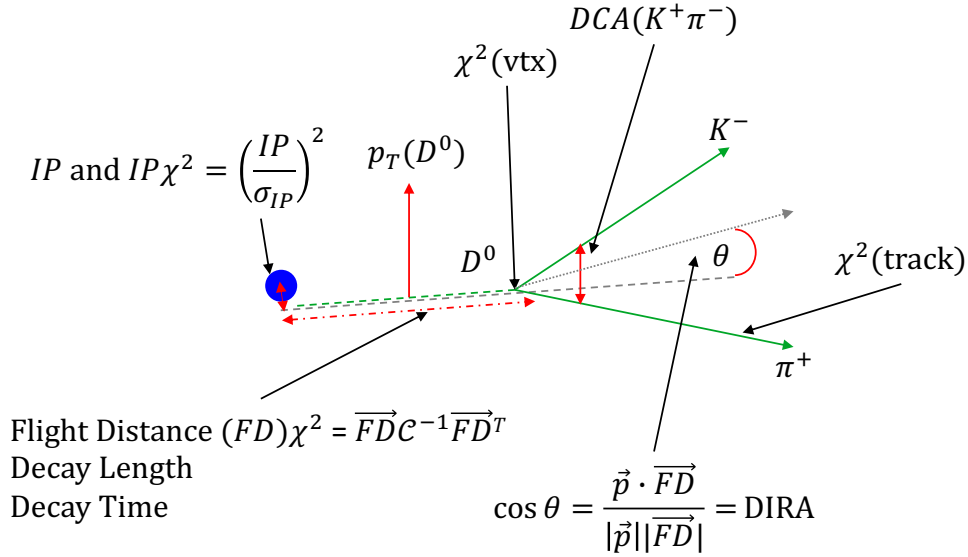


Figure 1: Visual description of the kinematic and topological variables used to select the  $D^0$  candidates.

183 trigger on this selection. For all events which passed this check, HFTRACKEFFICIENCY was  
 184 then run to match the true kaon and pion tracks to a reconstructed track within 5% for  $p_x$ ,  $p_y$   
 185 and  $p_z$  and write these tracks to a new track map. KFPARTICLE was then configured to run  
 186 on this new track map to ensure a pure  $D^0 \rightarrow K^-\pi^+$  sample. The only selection requirement  
 187 in KFPARTICLE was that the  $K^-\pi^+$  pair falls within the range  $1.70 \leq m_{K^-\pi^+} \leq 2.00$ .

188 The background sample was taken from a simulated sample of minimum bias Au+Au  
 189 collisions using HIJING. 10 thousand events were used with no selection other than the  
 190  $K^-\pi^+$  pair falls within the range  $2.00 \leq m_{K^-\pi^+} \leq 2.10$ .

191 ROOT's multi-variate analysis toolkit (TMVA) [19] can read in data sets where users  
 192 specify variables to use for multivariate analyses. Useful features of this kit is it can au-  
 193 tomatically overlay the signal and background distributions of each variable, measure the  
 194 variable correlations and show the agreement between training and testing samples. The sig-  
 195 nal and background samples were fed into TMVA where the event number that is stamped  
 196 by Fun4All is used to separate training and testing samples. The training samples used  
 197 exclusively odd numbered events while the testing samples used even numbered events. The  
 198 training and testing samples' invariant mass distributions are shown in Figure 2. 277005  
 199 candidates were used for the signal sample (138720 for training) and 884752 were used for  
 200 the background sample (449199 for training). The variable comparisons are shown in Fig-  
 201 ure 3 while the variable correlations are shown in Figure 4. TMVA was run only using  
 202 cuts-based methods, the associated receiver operating characteristic (ROC) plot is shown in  
 203 Figure 5. Using TMVA and assuming that there are 10000 background candidates to every  
 204 signal candidate, the package predictions an optimal signal efficiency of approximately 2%.

205 With this information, the baseline cuts to be applied to the minimum bias Au+Au  
 206 sample are given in Table 2. Before these cuts, there were 300755  $D^0$  signal candidates from  
 207 the  $p+p$  simulation and 1782139 background candidates from the Au+Au. After applying  
 208 these cuts, there were 39808  $D^0$  signal candidates from the  $p+p$  simulation and 978 back-



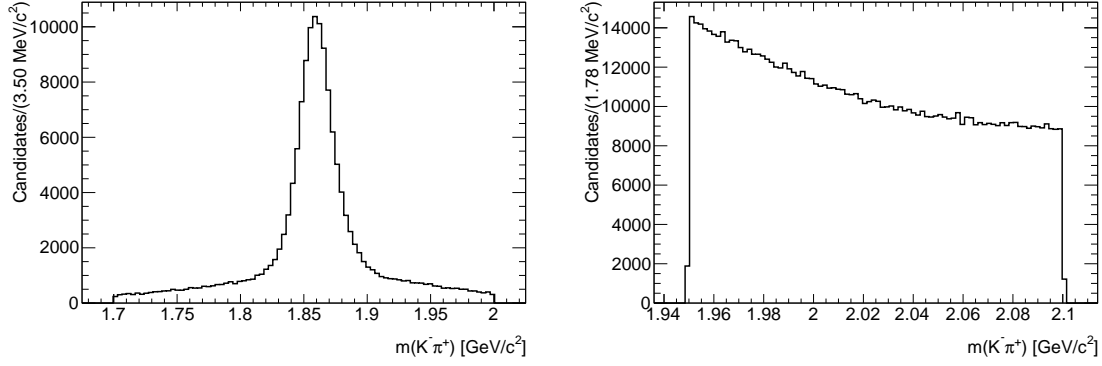


Figure 2: Invariant mass distributions of the signal and background samples used to make the selection decisions. The signal sample is shown on the left and the background sample is shown on the right.

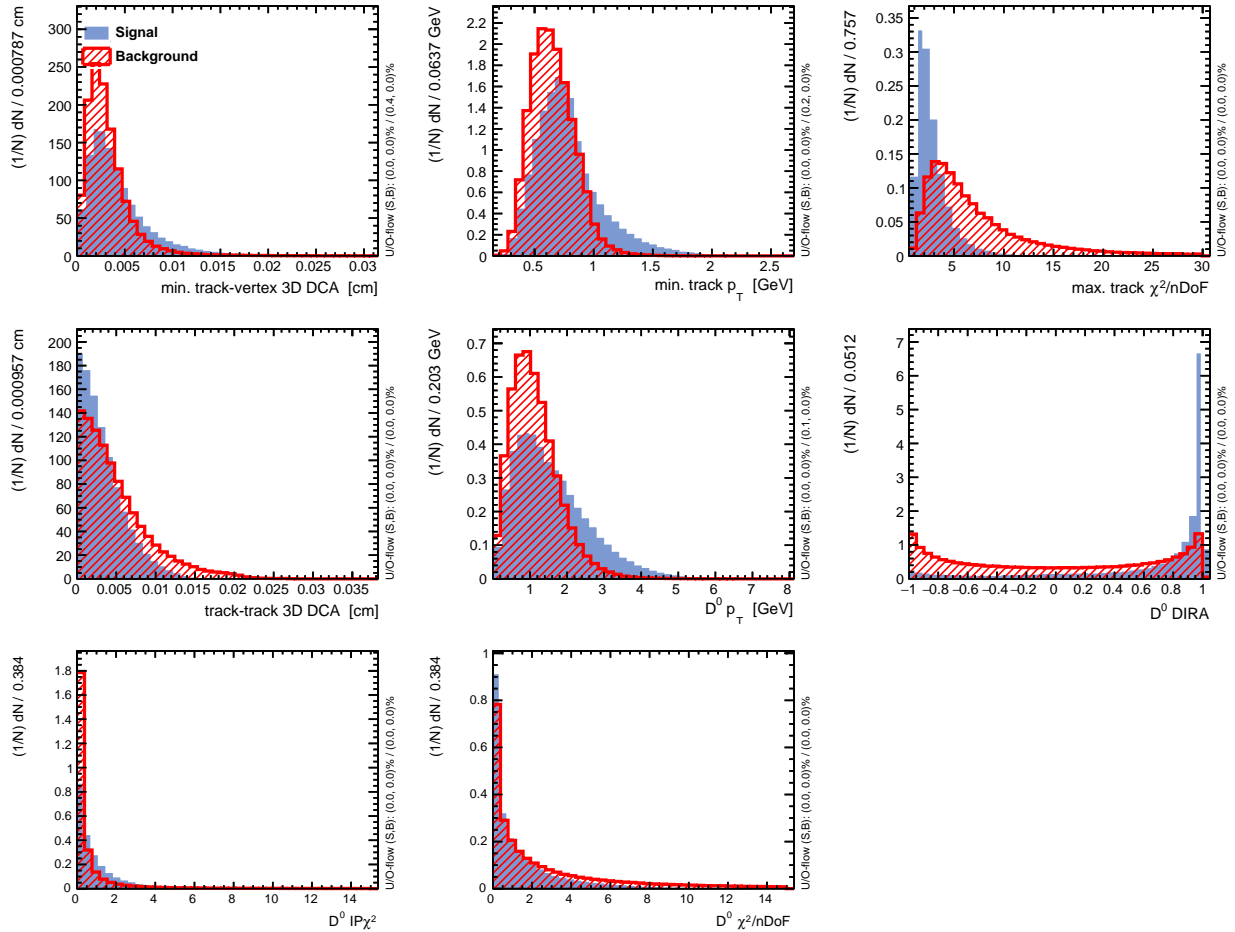


Figure 3: Input variable distributions of the signal (blue, whole) and background (red, dashed) samples used to make the selection decisions. Variables from top to bottom, left to right: minimum track IP, minimum track  $p_T$ , maximum track  $\chi^2$  per number of degrees of freedom, track-track DCA,  $D^0 p_T$ ,  $D^0$  DIRA,  $D^0$  IP  $\chi^2$ ,  $D^0 \chi^2$  per number of degrees of freedom.

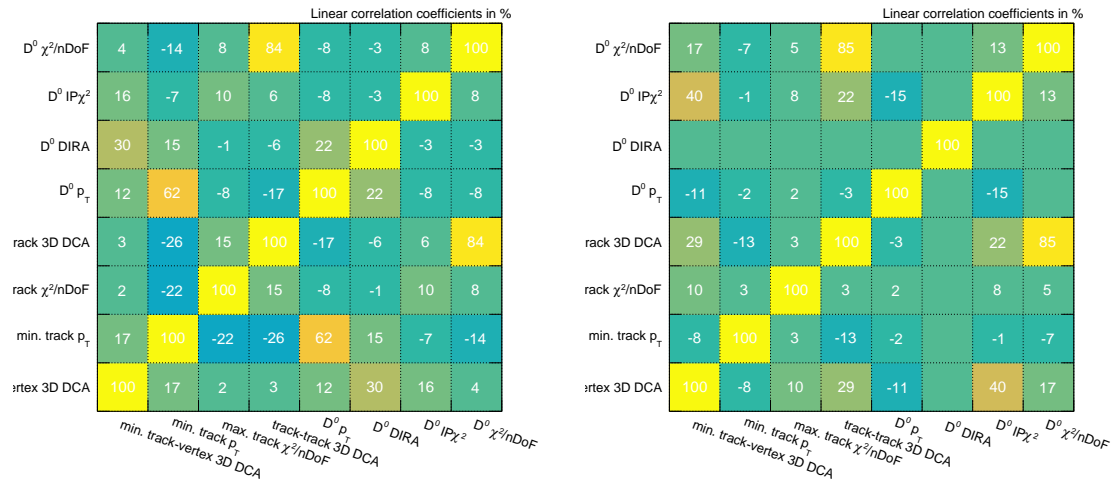


Figure 4: Input variable correlation coefficients of the signal (left) and background (right) samples used to make the selection decisions.

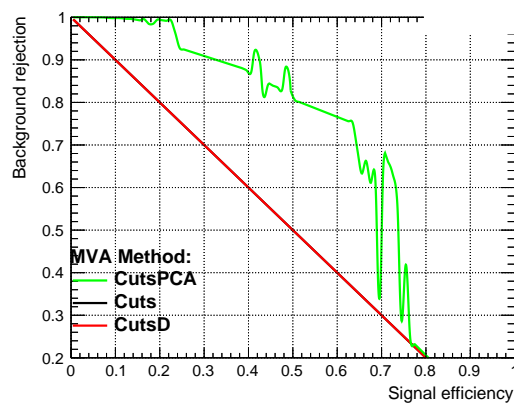


Figure 5: Receiver operating characteristic plot from the pre-selected variable distributions using cut methods.

209 ground candidates from the Au+Au. This gives a 13.24% efficiency on the signal and 0.05%  
 210 efficiency on the background selection. It should be noted that the invariant mass range in  
 211 Table 2 is slightly larger than what is used in the final fitting ( $1.70 \leq m_{K^- \pi^+} \leq 2.00$ ). This  
 212 is to obtain sufficient background samples on either side of the  $D^0$  mean mass to allow for a  
 213 machine-learning study.

Variable	Cut
min. track IP [cm]	0.0025
min. track $p_T$ [GeV]	0.7
max. track $\chi^2/\text{nDoF}$	5
max. track-track DCA [cm]	0.008
$D^0$ IP $\chi^2$	3
$D^0$ $p_T$ [GeV]	1.5
$D^0$ $\chi^2/\text{nDoF}$	5
$D^0$ DIRA	0.90
$m_{K^- \pi^+}$ [GeV]	1.65 $\rightarrow$ 2.10

Table 2: Baseline cuts used to select  $D^0 \rightarrow K^- \pi^+$  candidates.

## 214 5 Models

215 The modelling of the  $K^- \pi^+$  invariant mass distribution is performed using the RooFIT  
 216 fitting package [20] provide with ROOT.

### 217 5.1 Signal

218 As there is no PID at sPHENIX, it is possible to reconstruct a  $D^0$  but invert the mass  
 219 hypothesis for both tracks. The two hypotheses will return different values for the mass.  
 220 However, if the mass window is sufficiently large, it's increasingly likely to select the incorrect  
 221 combination. Due to this effect, the invariant mass of the  $K^- \pi^+$  pair will appear as the  
 222 composition of two Gaussian functions, one with a very large width. The PDF describing  
 223 the signal shape is

$$f(x; \mu, \sigma_{\text{cor-ID}}, \sigma_{\text{mis-ID}}, k_{\text{cor-ID}}) = N \cdot \left[ k_{\text{cor-ID}} \exp\left(-\frac{(x - \mu)^2}{2\sigma_{\text{cor-ID}}^2}\right) + (1 - k_{\text{cor-ID}}) \exp\left(-\frac{(x - \mu)^2}{2\sigma_{\text{mis-ID}}^2}\right) \right] \quad (3)$$

224 where  $\mu$  is the mean value,  $\sigma_{\text{cor-ID}}$  is the width of the Gaussian describing the candidates  
 225 with the correct mass hypothesis,  $\sigma_{\text{mis-ID}}$  is the width of the Gaussian describing the candi-  
 226 dates with the incorrect mass hypothesis and  $k_{\text{cor-ID}}$  is the fraction of candidates belonging  
 227 to the Gaussian with the correct mass hypothesis. When fitting the signal extracted from

228 the  $p+p$  baseline, all parameters are left floating. However, when fitting the final distribu-  
 229 tion  $k_{\text{cor-ID}}$  is fixed to the value from the  $p+p$  simulation and  $\sigma_{\text{mis-ID}}$  is scaled from  $\sigma_{\text{cor-ID}}$   
 230 by the ratio of widths as measured in the  $p+p$  simulation to account for the difference in  
 231 momentum resolution between  $p+p$  and Au+Au events

$$\sigma_{\text{mis-ID}}^{\text{Au+Au}} = \sigma_{\text{cor-ID}}^{\text{Au+Au}} \frac{\sigma_{\text{mis-ID}}^{p+p}}{\sigma_{\text{cor-ID}}^{p+p}} \quad (4)$$

232 The results of the fit to the signal extracted from  $p+p$  simulated events are given in  
 Table 3 and Figure 6 respectively.

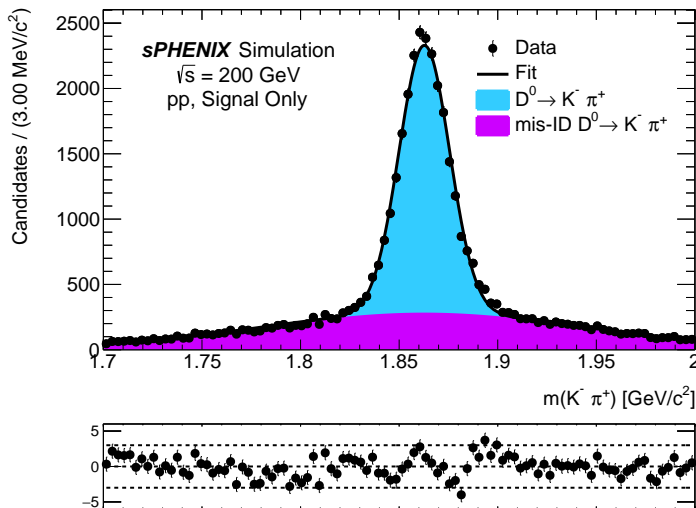


Figure 6: Default fit to the  $K^- \pi^+$  invariant mass distribution using simulated  $p+p$  events with truth matching.

233

Parameter	Value
$\mu$ [MeV]	$1862.78 \pm 0.11$
$\sigma_{\text{cor-ID}}$ [MeV]	$12.83 \pm 0.12$
$k_{\text{cor-ID}}$ [%]	$54.62 \pm 0.43$
$\sigma_{\text{mis-ID}}$ [MeV]	$82.89 \pm 1.00$

Table 3: Default fit parameters to the  $D^0 \rightarrow K^- \pi^+$   $p+p$  signal sample.

## 234 5.2 Background

235 The background is modeled using an exponential function as it was felt this would have a  
 236 better handle on the combinatorial background from the Au+Au simulation than a linear  
 237 function would. The background PDF is

$$f(x; k) = N \cdot \exp(\lambda x) \quad (5)$$

238 where  $\lambda$  is the decay constant.

### 239 5.2.1 Additional background sources

240 At this stage, there are no additional background sources modeled. However, there are  
241 several sources that could contribute:

- 242 •  $D^0 \rightarrow K^- \pi^+ \pi^0$  where the diphoton daughters are not used in the reconstruction. The  
243 branching ratio of  $D^0 \rightarrow K^- \pi^+ \pi^0$  to  $D^0 \rightarrow K^- \pi^+$  is approximately 3.5 and would lie  
244 to the lower mass region.
- 245 •  $D^+ \rightarrow K^- \pi^+ \pi^+$  where one of the pions are missed. The fragmentation fraction times  
246 branching ratio of  $D^+ \rightarrow K^- \pi^+ \pi^+$  to  $D^0 \rightarrow K^- \pi^+$  is approximately 1.0 and would  
247 lie to the lower mass region. The decay  $D^+ \rightarrow K^- K^+ \pi^+$ , where a kaon is missing is  
248 unlikely to contribute to the background in this study as either the mass loss from the  
249 kaon would be too great or the mis-ID of one kaon track as a pion as well as a missing  
250 track would likely make the channel fall out of the invariant mass window.
- 251 •  $D_s^+ \rightarrow K^- K^+ \pi^+$  where one of the kaons are missed. The fragmentation fraction times  
252 branching ratio of  $D_s^+ \rightarrow K^- K^+ \pi^+$  to  $D^0 \rightarrow K^- \pi^+$  is approximately 0.5 and would  
253 lie close to the  $D^0$  peak.
- 254 •  $\Lambda_c^+ \rightarrow p K^+ \pi^+$  is unlikely to contribute either due to the mis-ID of one track or the  
255 missing proton, which carries a significant portion of the energy.

256 The contribution of these background sources to this channel could be studied using  
257 dedicated simulations of these decays where the  $D^0 \rightarrow K^- \pi^+$  selection is applied and the  
258  $K^- \pi^+$  mass hypothesis is applied to the selection. The invariant mass distributions would  
259 not follow a predictable shape and so the Kernel method could be applied to model the  
260 shapes [21].

## 261 6 Results

262 The baseline selection in Table 2 was applied to the Au+Au simulation without using  
263 DECAFINDER and HFTRACKEFFICIENCY as a trigger to allow for a realistic background  
264 contamination. DECAFINDER and HFTRACKEFFICIENCY were both run alongside KFPARTICLE  
265 to calculate the acceptance, tracking, and selection efficiencies. The resulting  $m_{K^- \pi^+}$  spec-  
266 trum is shown in Figure 7 where there is no apparent  $D^0$  mass peak. The baseline cuts were  
267 chosen to be very loose so it was expected that they would need to be tightened at a later  
268 stage. The loose cuts were chosen to allow for a machine learning study which is detailed in  
269 Appendix B.

270 DECAFINDER and HFTRACKEFFICIENCY were run requiring that a maximum track  
271  $\eta \leq 1.6$  and minimum track  $p_T \geq 0.16$  GeV. HFTRACKEFFICIENCY also required that  
272 the true value of  $p_x$ ,  $p_y$  and  $p_z$  match to the reconstructed values within 5%. The  $p_T$  and  
273  $\eta$  values are right at the periphery of the sPHENIX tracking abilities and only significantly  
274 displaced tracks can meet these requirements or else they become loopers or fall outside of

275 the tracking acceptance. By using these values, the acceptance and tracking efficiency can be  
 276 underestimated by counting decays that could never be reconstructed. Thus, an improvement  
 277 in calculating the acceptance and tracking efficiency would be to add a recalculation of the  
 278  $p_T$  and  $\eta$  thresholds based on the position of the secondary vertex. This is outside of the  
 279 scope of this study but a second calculation of the tracking efficiency is performed with a  
 280 maximum track  $\eta \leq 0.5$  and minimum track  $p_T \geq 1.00$  GeV. Further, the track matching  
 281 requirement can reject lower  $p_T$  tracks where the momentum calculation is less accurate.

282 Using the loose requirements to accept a decay, the geometric acceptance is calculated  
 283 to be 34.5% and the tracking efficiency is calculated to be 7.8%. The geometric acceptance  
 284 includes  $D^0 \rightarrow K^- \pi^+$  decays that have an associated photon. The energy of these photons  
 285 is not known and so it is unknown if these decays could be meaningfully reconstructed if the  
 286 energy loss is too great. If all decays with an associated photon are assumed to fall outside  
 287 of the invariant mass window, then the acceptance is reduced to 24.3%. Table 4 details  
 288 the number of generated decays and how they do or do not fall inside the loose sPHENIX  
 289 acceptance. If the tracking requirement is tightened to the values in the previous paragraph,  
 290 then the efficiency increases by a factor of 3 to 20.4%. Table 5 also details the number of  
 291 decays with an associated photon within the loose cut and how many had all their tracks  
 292 reconstructed via the tracking algorithms. A comparison of the tracking efficiency as a 2D  
 293 plot of minimum and maximum track  $p_T$  and the efficiency as a function of the true  $D^0$   $p_T$   
 294 is given in Figure 8.

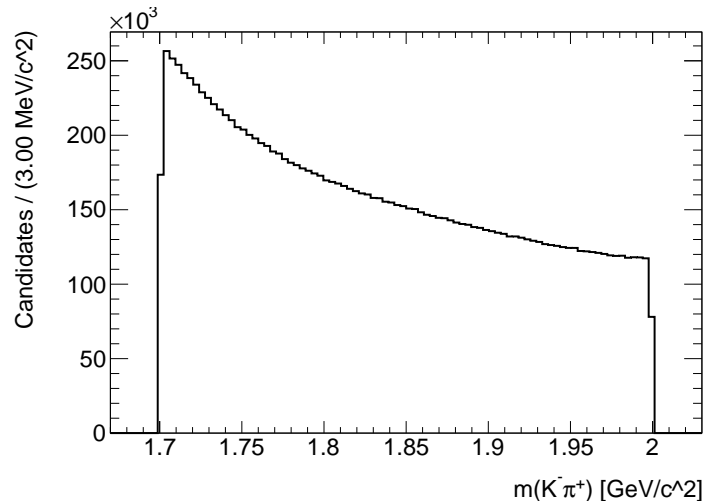


Figure 7: Distribution of the  $K^- \pi^+$  invariant mass pairs using simulated Au+Au events with the baseline selection.

295 A large volume of the  $D^0$  daughter tracks appear to sit at the edge or beyond the  
 296 sPHENIX fiducial region where one track has a either large value of  $|\eta|$  or a low  $p_T$  (or  
 297 both). This can be seen in Figure 9 where there is a large concentration of events with a  $p_T$   
 298 less than 0.2 GeV and an  $|\eta| \geq 1.1$ .

299 A tighter selection was applied to the data set to extract the  $D^0$  peak. The cuts were  
 300 chosen by studying the kinematic distributions of signal and background samples after the  
 301 baseline selection was applied. To avoid biases, the  $p+p$  sample was used for the signal

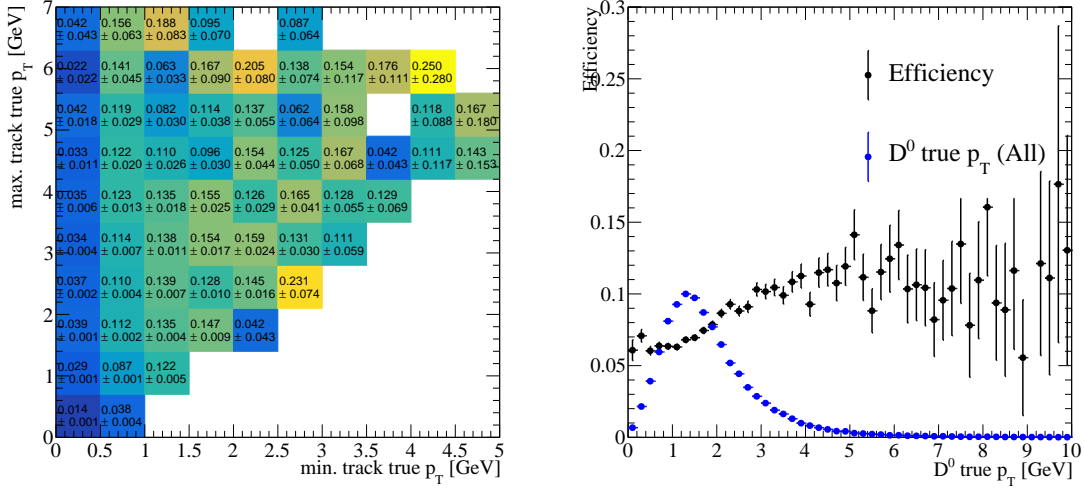


Figure 8: The tracking efficiency as measured in the Au+Au simulation. A 2D plot of the tracking efficiency as a function of the minimum track  $p_T$  (x-axis) and maximum track  $p_T$  (y-axis) is shown on the left while the tracking efficiency as a function of the true  $D^0$   $p_T$  is shown on the right..

Requirement	Count
Generated decays	636070
Generated decays that fail $p_T$	16766
Generated decays that fail $\eta$	300839
Generated decays that fail $p_T$ and $\eta$	98869
Reconstructable decays	219596
Reconstructable decays with an associated $\gamma$	65048
Reconstructable decays with an associated $\pi^0$	0
Reconstructable decays with an associated $\gamma$ and $\pi^0$	0
Reconstructable decays with no associated $\gamma$ no $\pi^0$	154548
Efficiency [%]	34.5

Table 4: Comparison of the generated number of  $D^0 \rightarrow K^- \pi^+$  decays in the Au+Au simulation to the number of decays that fall in the sPHENIX acceptance of  $|\eta| \leq 1.6$  and  $p_T \geq 0.16$  MeV.

Requirement	Count
Decays in acceptance	219596
Decays fully reconstructed	17138
Efficiency [%]	7.8

Table 5: Comparison of the number of  $D^0 \rightarrow K^- \pi^+$  decays fully in the sPHENIX acceptance in the Au+Au simulation to the number of decays that have all track reconstructed.

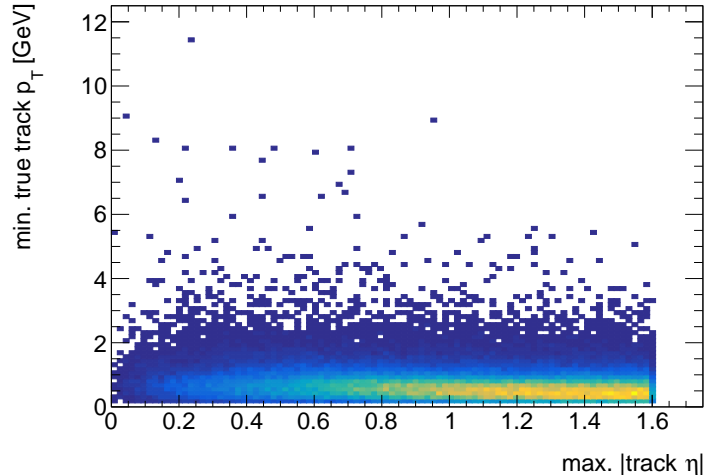


Figure 9: 2D distribution of the maximum absolute  $\eta$  value of a track and the minimum  $p_T$  of a track. Both of these values come from the truth information and do not necessarily come from the same track. Only one track needs to fall outside the fiducial range to be unable to reconstruct a decay.

302 and the upper and lower mass sidebands from the Au+Au sample that sit outside of the fit  
 303 range were used for the background sample. These samples and the corresponding variable  
 304 distributions can be found in Appendix B while the new cuts are listed in Table 6.

305 Two fits were performed using the models detailed in Section 5 at two different values  
 306 of the  $D^0$   $p_T$ : 2 and 4 GeV. The latter cut is determined to be a very tight cut capable of  
 307 rejecting a significant fraction of the background and could be achieved at an early stage  
 308 of the commissioning period while the tracking calibrations are being studied. The former  
 309 cut represents a physics goal of sPHENIX, to reconstruct  $D^0$  to a  $p_T$  of 2 GeV or less. The  
 310 results of the fit for  $D^0$   $p_T \geq 4$  GeV are given in Table 7 and Figure 10 where the yield was  
 311 measured to be  $129 \pm 20$ . This is a statistical significance of  $\sigma = 6.45$ . The results of the  
 312 fit for  $D^0$   $p_T \geq 2$  GeV are given in Table 8 and Figure 11 where the yield was measured to  
 313 be  $536 \pm 79$ . This is a statistical significance of  $\sigma = 6.78$ .

Variable	Cut
min. track IP [cm]	0.008
max. track-track DCA [cm]	0.005
$D^0$ $p_T$ [GeV]	2 or 4
$D^0$ DIRA	0.98
$m_{K^-\pi^+}$ [GeV]	1.70 $\rightarrow$ 2.00

Table 6: Enhanced cuts used to select  $D^0 \rightarrow K^-\pi^+$  candidates.



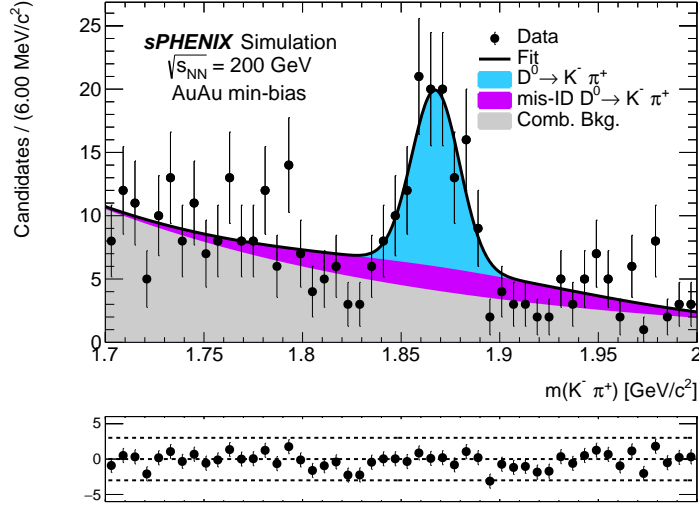


Figure 10: Fit to the  $K^-\pi^+$  invariant mass distribution using simulated Au+Au events with the tighter selection and a  $D^0 p_T \geq 4$  GeV.

Parameter	Value
$\mu$ [MeV]	$1867.45 \pm 2.1$
$\sigma_{\text{cor-ID}}$ [MeV]	$12.1 \pm 1.7$
$\lambda$ [ $\text{MeV}^{-1}$ ]	$-5.6 \pm 1.0$
$N_{\text{cand}}$	383
$f_{D^0}$ [%]	$33.7 \pm 5.3$

Table 7: Fit results with the increased selection requirements and a  $D^0 p_T \geq 4$  GeV.

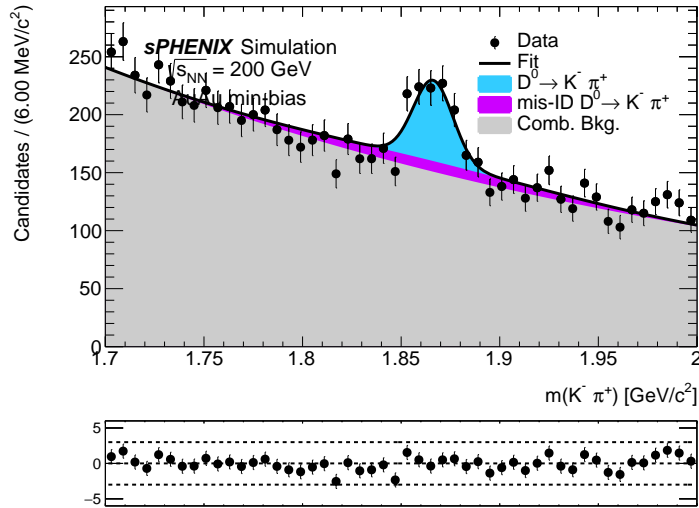


Figure 11: Fit to the  $K^-\pi^+$  invariant mass distribution using simulated Au+Au events with the tighter selection and a  $D^0 p_T \geq 2$  GeV.

Parameter	Value
$\mu$ [MeV]	$1866.23 \pm 1.6$
$\sigma_{\text{cor-ID}}$ [MeV]	$10.1 \pm 1.5$
$\lambda$ [MeV <sup>-1</sup> ]	$-2.80 \pm 0.14$
$N_{\text{cand}}$	8664
$f_{D^0}$ [%]	$6.2 \pm 0.9$

Table 8: Fit results with the increased selection requirements and a  $D^0$   $p_T \geq 2$  GeV.

## 7 Conclusion

This note details a study performed using simulated Au+Au events with the sPHENIX detector, and tracking and reconstruction algorithms that will be deployed during the commissioning period which is due to commence in Spring 2023. Approximately 22 million minimum-bias events were simulated, corresponding to just over 20 minutes of data taking at the full RHIC collision rate of 15 kHz. This study is aimed at defining base cuts to be deployed during the commissioning period where the collision rate is lower and so this sample corresponds to a longer integrated time than 20 minutes.

By applying the cuts detailed in Tables 2 and 6 and the models described in Section 5,  $D^0 \rightarrow K^-\pi^+$  decays were extracted with a significance larger than  $5\sigma$  for a mother  $p_T \geq 2$  GeV. It is currently recommended that these cuts be applied to the initial sPHENIX data sample to obtain a sufficient quantity and quality of  $D^0$  candidates to validate the tracking and heavy flavor programs for full physics analyses.

There are several areas where this study could be improved and some of these are detailed in the appendices in this note such as the application of machine learning, improved background and signal modelling, and the use of statistical techniques to unfold data for reaching a lower  $p_T$  region than 2 GeV. There are other areas where this study could be improved that are not discussed in detail of this note. One area is the improvement of the acceptance and tracking efficiency calculations. It was mentioned in Section 6 that the efficiency calculations do not account for the position of the secondary vertex and how this can alter the track  $\eta$  and  $p_T$  requirements. A module that can calculate the position dependent requirements would improve the understanding of these efficiencies. Further, the calculated values of the DCA variables appears to be smaller than would be expected for a particle that travels around 1 mm before decaying. A study will be performed that compares the DCA values calculated by KFPARTICLE with a typical straight line extrapolation of a track to the primary vertex to see if there are any discrepancies.

Beyond the additional studies detailed in the appendices, this work could be expanded easily to study the  $K^+K^-\pi^+$  spectrum to obtain  $D_{(s)}^+$  candidates with little alterations to the selection beyond the invariant mass window. This study could also be used to serve as a baseline for a  $\Lambda_c^+$  study (which has a shorter lifetime than the  $D^0$ ) and a  $K_s^0$  study (which has a longer lifetime than the  $D^0$  but a lower  $p_T$  spectrum).

## References

- [1] sPHENIX Collaboration. sPHENIX Beam Use Proposal, May 2022.
- [2] G. Aad et al. Measurement of the nuclear modification factor for muons from charm and bottom hadrons in Pb+Pb collisions at 5.02 TeV with the ATLAS detector. *Physics Letters B*, 829:137077, 2022.
- [3] ALICE Collaboration. Measurement of beauty production via non-prompt  $D^0$  mesons in Pb-Pb collisions at  $\sqrt{s_{NN}} = 5.02$  TeV, 2022.
- [4] X. Chen. Prompt and non-prompt  $D^0$ -meson production in Au+Au collisions, Ph.D. Thesis, STAR, 2019.
- [5] P. Jackson et al. Measurement of the total cross section from elastic scattering in pp collisions at  $\sqrt{s} = 8$  TeV with the ATLAS detector. *Physics Letters. Section B: Nuclear, Elementary Particle and High-Energy Physics*, 761:158–178, 2016.
- [6] R. Aaij et al. Prompt charm production in pp collisions at  $\sqrt{s} = 7$  TeV. *Nuclear Physics B*, 871(1):1–20, 2013.
- [7] S. Acharya et al. Charm-quark fragmentation fractions and production cross section at midrapidity in pp collisions at the LHC. *Physical Review D*, 105(1), jan 2022.
- [8] Particle Data Group. Review of Particle Physics. *Progress of Theoretical and Experimental Physics*, 2020(8), 08 2020. 083C01.
- [9] T. Sjöstrand, S. Mrenna, and P. Skands. A Brief Introduction to PYTHIA 8.1. *Computer Physics Communications*, 178(11):852–867, 2008.
- [10] S. Lim, W Park. PYTHIA8 tune in pp 200 GeV. Presentation at Heavy Flavor Topical Group Meeting, Dec. 14, 2020, [https://indico.bnl.gov/event/10309/contributions/44139/attachments/31909/50542/sPHENIX\\_HF\\_shlim\\_20201215.pdf](https://indico.bnl.gov/event/10309/contributions/44139/attachments/31909/50542/sPHENIX_HF_shlim_20201215.pdf).
- [11] M. Gyulassy and X.-N. Wang. HIJING 1.0: A Monte Carlo program for parton and particle production in high energy hadronic and nuclear collisions. *Computer Physics Communications*, 83(2):307–331, 1994.
- [12] S. Agostinelli et al. Geant4—a simulation toolkit. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*, 506(3):250–303, 2003.
- [13] X. Ai et al. A Common Tracking Software Project, 2021.
- [14] J. D. Osborn, A. D. Frawley, J. Huang, S. Lee, H. Pereira Da Costa, M. Peters, C. Pinkenburg, C. Roland, and H. Yu. Implementation of ACTS into sPHENIX Track Reconstruction. *Computing and Software for Big Science*, 5(1), Oct 2021.

- 379 [15] A. Gorbunov and I. Kisel. Reconstruction of decayed particles based on the Kalman  
380 filter. *CBM-SOFT-note-2007-003*, May 2007.
- 381 [16] C. Dean. Event Reconstruction and MDC1. Presentation at sPHENIX 10<sup>th</sup>  
382 Collaboration Meeting, Jan. 22, 2021, [https://indico.bnl.gov/event/10568/  
383 contributions/45092/attachments/32421/51581/HeavyFlavor\\_EventRecoAndMDC\\_  
384 C\\_Dean\\_20210122.pdf](https://indico.bnl.gov/event/10568/contributions/45092/attachments/32421/51581/HeavyFlavor_EventRecoAndMDC_C_Dean_20210122.pdf).
- 385 [17] S. T. Araya, C. Dean, J. Huang, H. Okawa, , and Z. Shi. First MDC1 Results from  
386 Heavy Flavor Topical Group, April 2021.
- 387 [18] A. Buckley, P. Ilten, D. Konstantinov, L. Lönnblad, J. Monk, W. Pokorski, T. Przedzin-  
388 ski, and A. Verbytskyi. The HepMC3 event record library for Monte Carlo event gen-  
389 erators. *Computer Physics Communications*, 260:107310, March 2021.
- 390 [19] A. Hoecker et al. TMVA - Toolkit for Multivariate Data Analysis, 2007.
- 391 [20] W Verkerke and D. P. Kirkby. The RooFit toolkit for data modeling. *eConf*,  
392 C0303241:MOLT007, 2003. [arXiv:physics/0306116](https://arxiv.org/abs/physics/0306116).
- 393 [21] K. Cranmer. Kernel Estimation in High-Energy Physics. *Computer Physics Communi-  
394 cations*, 136(3):198–207, 2001.
- 395 [22] J. Gaiser. Charmonium Spectroscopy From Radiative Decays of the  $J/\psi$  and  $\psi'$ . *Ph.D.  
396 Thesis, SLAC*, 1982.
- 397 [23] B. Knuteson and H. Miettinen. Mass Analysis and Parameter Estimation with PDE.  
398 *D0 notes*, 9 1997.
- 399 [24] I. S. Abramson. On Bandwidth Variation in Kernel Estimates - A Square Root Law.  
400 *The Annals of Statistics*, pages 1217–1223, 1982.
- 401 [25] M. Pivk and F.R. Le Diberder. sPlot: A Statistical Tool to Unfold Data Distributions.  
402 *Nuclear Instruments and Methods in Physics Research A*, 555(1):356 – 369, 2005.

# 403 Appendices

## 404 A Alternative Fit Models

405 In Figure 6 it can be seen that the default model does not perfectly describe the data as  
 406 it overshoots the mean while the base seems to be shifted in the opposite direction (it over  
 407 fits before the mean and under fits after the mean). This implies that simply adjusting the  
 408 fit parameters will not help. Alternative fit models are proposed that can better describe  
 409 the data. The fits are described only to emphasise the default fit shortcomings. The double  
 410 Gaussian model is used in the final fit as the  $D^0$  shape is known to be well described by  
 411 a single Gaussian. The reason for the asymmetric peak shape is assumed to be due to an  
 412 earlier version of the tracking that was run on the  $p+p$  simulation compared to the Au+Au  
 413 simulation. This earlier version of the tracking had a simpler version of the TPC clustering  
 414 which was improved for the Au+Au simulation and was known to affect the momentum  
 415 resolution of the tracks.

### 416 A.1 Bifurcated Gaussian

417 The first model used is a bifurcated Gaussian which has different widths on either side of  
 418 the mean value. This shape is described by the PDF

$$f(x; \mu, \sigma_L, \sigma_R) = N \cdot \begin{cases} \exp\left(-\frac{(x - \mu)^2}{2\sigma_L^2}\right), & \text{for } x < \mu \\ \exp\left(-\frac{(x - \mu)^2}{2\sigma_R^2}\right), & \text{for } x \geq \mu \end{cases} \quad (6)$$

419 where  $\mu$  describes the mean of the Gaussian and  $\sigma_{[L/R]}$  describe the Gaussian widths on the  
 420 left and right hand side of the mean value.

421 The model is fit to a subset of the  $p+p$  signal sample. This subset has no cuts applied  
 422 and uses 38420 signal candidates in the fit. Using a smaller number of events avoids the issue  
 423 of over-fitting the distribution but having no selection means the shape will differ from the  
 424 final distribution. The fit values and plot are given in Table 9 and Figure 12 respectively. It  
 425 can be seen that the alternative fit model better describes the data, however there is still a  
 426 prominent bias.

Parameter	Value
$\mu$ [MeV]	$1858.93 \pm 0.28$
$\sigma_{\text{cor-ID,L}}$ [MeV]	$12.02 \pm 0.22$
$\sigma_{\text{cor-ID,R}}$ [MeV]	$14.98 \pm 0.25$
$k_{\text{cor-ID}}$ [%]	$57.73 \pm 0.44$
$\sigma_{\text{mis-ID}}$ [MeV]	$89.21 \pm 1.31$

Table 9: Bifurcated Gaussian fit parameters to the  $D^0 \rightarrow K^-\pi^+$   $p+p$  signal sample.

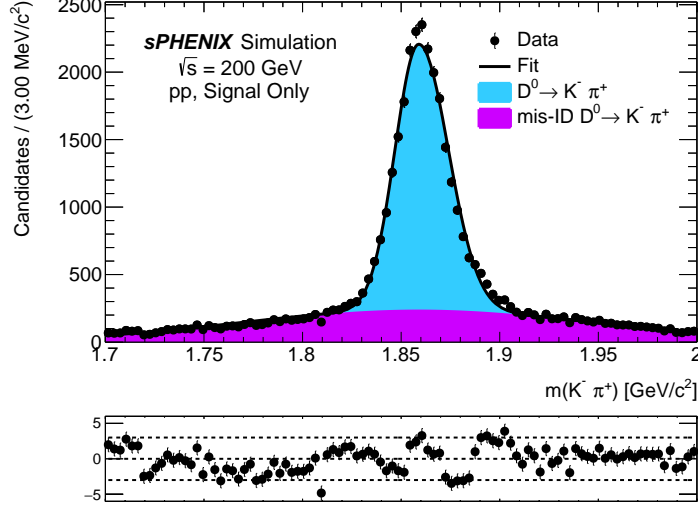


Figure 12: Bifurcated Gaussian fit to the  $K^- \pi^+$  invariant mass distribution using simulated  $p+p$  events with truth matching.

## 427 A.2 Double Crystal Ball

428 A second alternative model is proposed to describe the data. This model is a double-sided  
 429 Crystal Ball function [22] and consists of a Gaussian core with a polynomial tail. This  
 430 tail is capable of describing radiative losses. The double-sided Crystal Ball has this tail on  
 431 either side of the Gaussian core unlike the single-sided which only has the polynomial on the  
 432 low-mass mass side. The shape is described by the PDF

$$f(x; \alpha_L, n_L, \alpha_H, n_H, \mu, \sigma) = N \cdot \begin{cases} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), & \text{for } -\alpha_L < \frac{x-\mu}{\sigma} < -\alpha_H \\ A_L \cdot \left(B_L - \frac{x-\mu}{\sigma}\right)^{-n_L}, & \text{for } -\alpha_L \geq \frac{x-\mu}{\sigma} \\ A_H \cdot \left(B_H - \frac{x-\mu}{\sigma}\right)^{-n_H}, & \text{for } \frac{x-\mu}{\sigma} \geq -\alpha_H \end{cases} \quad (7)$$

433 where

$$A_{[L/H]} = \left( \frac{n_{[L/H]}}{|\alpha_{[L/H]}|} \right)^{n_{[L/H]}} \cdot \exp\left(-\frac{|\alpha_{[L/H]}|^2}{2}\right) \quad (8)$$

$$B_{[L/H]} = \frac{n_{[L/H]}}{|\alpha_{[L/H]}|} - |\alpha_{[L/H]}|. \quad (9)$$

434 The mean and width of the central Gaussian distribution are given by  $\mu$  and  $\sigma$  respec-  
 435 tively while L and H refer to the parameters defining the low and high mass regions of the  
 436 distribution. The parameter  $\alpha_{[L/H]}$  describes the boundary between the Gaussian and power  
 437 law components while  $n_{[L/H]}$  describes the order of the power law (it is not necessarily an  
 438 integer value).

439 The model is fit to the same subset of the  $p+p$  signal sample. The fit values and plot are  
 440 given in Table 10 and Figure 13 respectively. It can be seen that the double Crystal Ball  
 441 describes the data well.

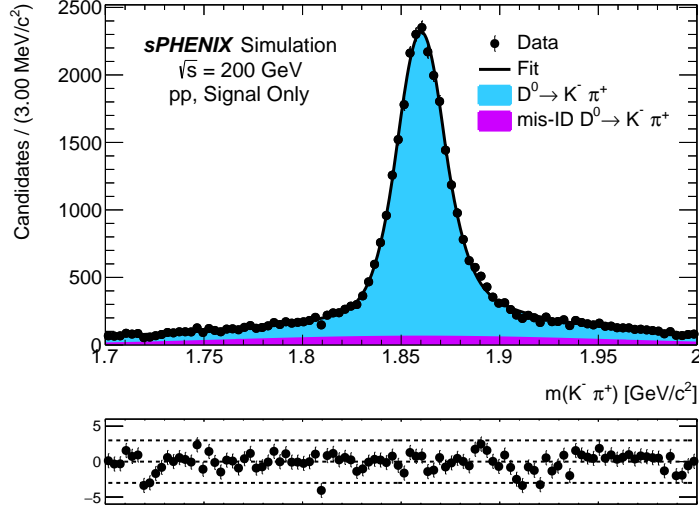


Figure 13: Double Crystal Ball fit to the  $K^- \pi^+$  invariant mass distribution using simulated  $p+p$  events with truth matching.

Parameter	Value
$\mu$ [MeV]	$1860.15 \pm 0.12$
$\sigma_{\text{cor-ID}}$ [MeV]	$12.67 \pm 0.15$
$\alpha_L$	$1.07 \pm 0.03$
$n_L$	$0.89 \pm 0.20$
$\alpha_H$	$0.93 \pm 0.03$
$n_H$	$1.31 \pm 0.30$
$k_{\text{DCB, Gauss. core}}$ [%]	$41.53 \pm 2.86$
$k_{\text{cor-ID}}$ [%]	$87.29 \pm 6.26$
$\sigma_{\text{mis-ID}}$ [MeV]	$97.7 \pm 19.4$

Table 10: Double Crystal Ball fit parameters to the  $D^0 \rightarrow K^- \pi^+$   $p+p$  signal sample.

### 442 A.3 Kernel density estimated backgrounds

443 No partially reconstructed backgrounds are modelled in the default fit. However, if one  
 444 wishes to model these non-parametric distributions an option would be to use a kernel  
 445 density estimation (KDE). As the invariant mass distribution of the events that pass the  
 446 selection requirements tend to become non-parametric due to the selection requirements and  
 447 the invariant mass substitutions imposed on the sample, then the shape of the PDFs are  
 448 extracted using the kernel density estimation method with an adaptive bandwidth [21]. The  
 449 use of kernel density estimation to describe non-parametric distributions in invariant masses  
 450 was proposed by the D0 collaboration for analysis of the Higgs boson [23] but has been  
 451 developed for use in other collaborations.

452 In a kernel estimation the events of a distribution are substituted for a kernel function  
 453 so that the distribution,  $\hat{f}_0(x)$ , can be described as

$$\hat{f}_0(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-t_i}{h}\right) \quad (10)$$

454 where  $t_i$  is the value of event  $i$  and  $h$  is the bandwidth or smoothing parameter. It has  
 455 been suggested [23] that a suitable kernel function for use in describing the invariant mass  
 456 distribution would be a Gaussian as it is positive definite and infinitely differentiable.

457 To avoid issues of function overspill at boundaries, underestimations of the distribution  
 458 within regions with low event density, and overestimation within regions of high density then  
 459 the bandwidth is allowed to alter on an event-by-event basis which is known as an *adaptive*  
 460 *kernel estimation* where the per-event bandwidth is given by [24]

$$h_i = \frac{h}{\sqrt{\hat{f}_0(x)}}. \quad (11)$$

461 The PDFs determined by the kernel method were obtained by using DECAYFINDER and  
 462 HFTRACKEFFICIENCY to select  $D_s^+ \rightarrow K^- K^+ \pi^+$  from the  $p+p$  to  $c\bar{c}$  simulation. The new  
 463 track map of these decays was then passed to KFPARTICLE where the default selection  
 464 given in Table 6 was applied using the  $D^0 \rightarrow K^- \pi^+$  decay descriptor. The resulting shape is  
 465 given in Figure 14. The expected contamination of this decay with respect to  $D^0 \rightarrow K^- \pi^+$   
 466 would be estimated from the ratio of the hadronisation fractions, branching fractions and  
 467 the number of each decay that pass the selection divided by the total number of each decay  
 468 generated,

$$k_{\text{chann}}^{\text{mid}} = \frac{f_{\text{mid}}}{f_{\text{chann}}} \frac{\mathcal{B}_{\text{mid}}}{\mathcal{B}_{\text{chann}}} \frac{\omega_{\text{mid}}}{\varepsilon_{\text{chann}}} \quad (12)$$

469 where  $f_{\text{P}}$  is the hadronization fraction of the mother particle,  $\mathcal{B}_{\text{P}}$  is the branching fraction  
 470 of the decay,  $\omega_{\text{mid}}$  is the misidentification efficiency and  $\varepsilon_{\text{chann}}$  is the signal efficiency. As both  
 471 the signal and background samples are drawn from the same minimum bias  $p+p$  simulation  
 472 where all  $c\bar{c}$  events are saved, then the contamination fraction becomes

$$k_{\text{chann}}^{\text{mid}} = \frac{N_{\text{sel}}(D_s^+ \rightarrow K^- K^+ \pi^+)}{N_{\text{sel}}(D^0 \rightarrow K^- \pi^+)} \quad (13)$$



where  $N_{\text{sel}}(X)$  is the number of decays of  $X$  that pass the selection.

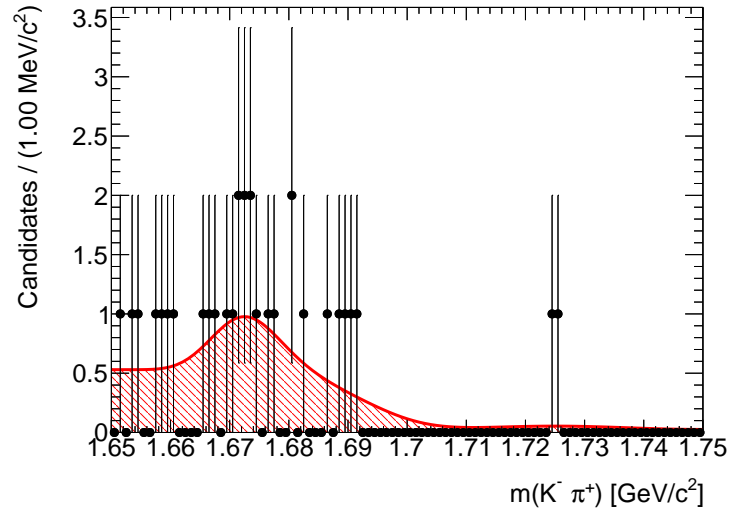


Figure 14: Kernel estimated distribution of  $D_s^+ \rightarrow K^- K^+ \pi^+$  in the  $m_{K^- \pi^+}$  spectrum using the tighter  $D^0 \rightarrow K^- \pi^+$  selection.

## 474 B Machine Learning Selection

475 When selecting the baseline cuts, it was known that they would be so loose that the back-  
 476 ground would overwhelm the signal. The reason behind this cut was to allow for a machine-  
 477 learning study to see if the low- $p_T$  reach of  $D^0$  could be improved. To this end, signal and  
 478 background samples were passed to ROOT's TMVA machine learning package [19] to train  
 479 various algorithms and see if this goal could be achieved. The signal sample was comprised  
 480 of the  $p+p$  sample with the baseline cuts applied while the background sample was taken  
 481 from the upper and lower mass side bands of the Au+Au simulation with the baseline se-  
 482 lection applied. The background sample was taken from the ranges  $1.65 \leq m_{K-\pi^+} < 1.70$   
 483 and  $2.00 < m_{K-\pi^+} \leq 2.10$ . The samples were split into even and odd samples by their event  
 484 numbers. The odd samples were used for training while the even samples were used for test-  
 485 ing. All of this was to avoid biases in the final selection by not re-using events. The number  
 486 of events used in each sample is given in Table 11 and the invariant mass distributions are  
 487 given in Figure 15.

Sample	Number of candidates
Signal, training	20082
Signal, testing	19726
Background, training	199228
Background, testing	200662

Table 11: Number of candidates used for training and testing the machine learning algorithms.

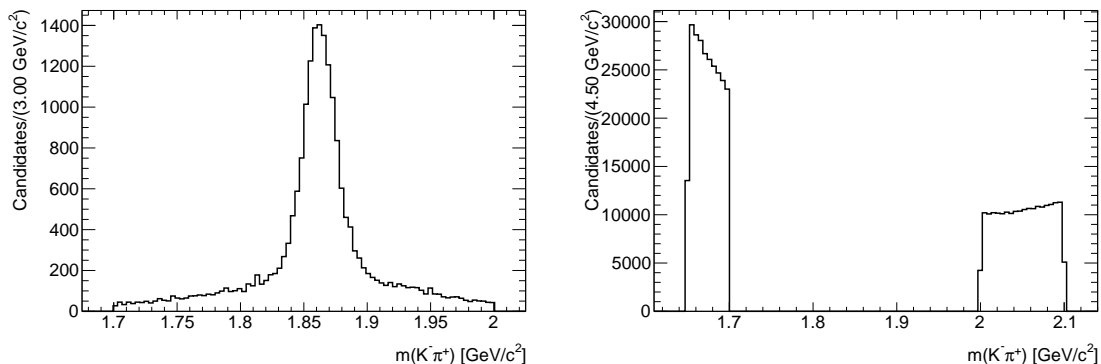


Figure 15: Invariant mass distributions of the signal and background samples used to train the machine learning algorithms. The preselection given in Table 2 has been applied to both samples. Each sample consists of events with an odd event number. The signal sample is shown on the left and the background sample is shown on the right.

488 Thirteen variables were chosen to train the algorithms: the minimum and maximum DCA  
 489 of the selected tracks with respect to their primary vertex in both two and three dimensions,  
 490 the minimum and maximum track  $p_T$ , the minimum and maximum track  $\chi^2$  per number of

491 degrees of freedom, the DCA of the tracks with respect to each other, the  $D^0$   $p_T$ , the  $D^0$   
492 DIRA, the  $D^0$  DCA  $\chi^2$  with respect to the primary vertex and the  $D^0$   $\chi^2$  per number of  
493 degrees of freedom. The comparison of the signal and background distributions for these  
494 variables can be seen in Figure 16 and the correlations between variables can be seen in  
495 Figure 17. These variables were used to train six algorithms: two multi-layer perceptrons  
496 (MLP), ROOTs own neural net, a boosted decision tree (BDT), a BDT with gradient boost  
497 and a BDT with de-correlation and adaptive boost. The receiver operating characteristic  
498 (ROC) curve and response of each algorithm with training and testing samples overlaid can  
499 be seen in Figures 18 and 19 respectively.

500 The ROC curve legend lists the algorithms in order of best performance. Thus, the old  
501 MLP is predicted to give the best performance. The algorithm was applied to the Au+Au  
502 data within the range  $1.70 \leq m_{K-\pi^+} \leq 2.00$  and the response of each candidate to each  
503 ML algorithm was calculated. The fit to the invariant mass was then performed using the  
504 baseline cut and two values of the MLP. For an MLP response  $\geq 0.99$ , the yield was measured  
505 to be  $536 \pm 74$  while for an MLP response  $\geq 0.999$ , the yield was measured to be  $339 \pm 48$ .  
506 The fit to both distributions is given in Figure 20 and the fit results with an algorithm  
507 response  $\geq 0.999$  are given in Table 12.

Parameter	Value
$\mu$ [MeV]	$1865.55 \pm 1.5$
$\sigma_{\text{cor-ID}}$ [MeV]	$10.4 \pm 1.4$
$\lambda$ [MeV $^{-1}$ ]	$-6.60 \pm 0.29$
$N_{\text{cand}}$	3139
$f_{D^0}$ [%]	$10.8 \pm 1.5$

Table 12: Fit results for the ML study with a multi-layer perceptron with an algorithm response  $\geq 0.999$ .

508 While this study was brief, it demonstrates that machine learning algorithms can be  
509 used to extract the  $D^0 \rightarrow K^-\pi^+$  decay at sPHENIX. This could be applied to the entire  
510 sPHENIX data set but it is not felt to be suited to the commissioning period where we wish  
511 to pre-select a small subset of tracks that we can quickly reprocess at a later date. The  
512 baseline selection resulted in 15 million  $D^0$  candidates from a sample of 21 million minimum  
513 bias events and thus would result in an unnecessarily large data sample to reprocess.

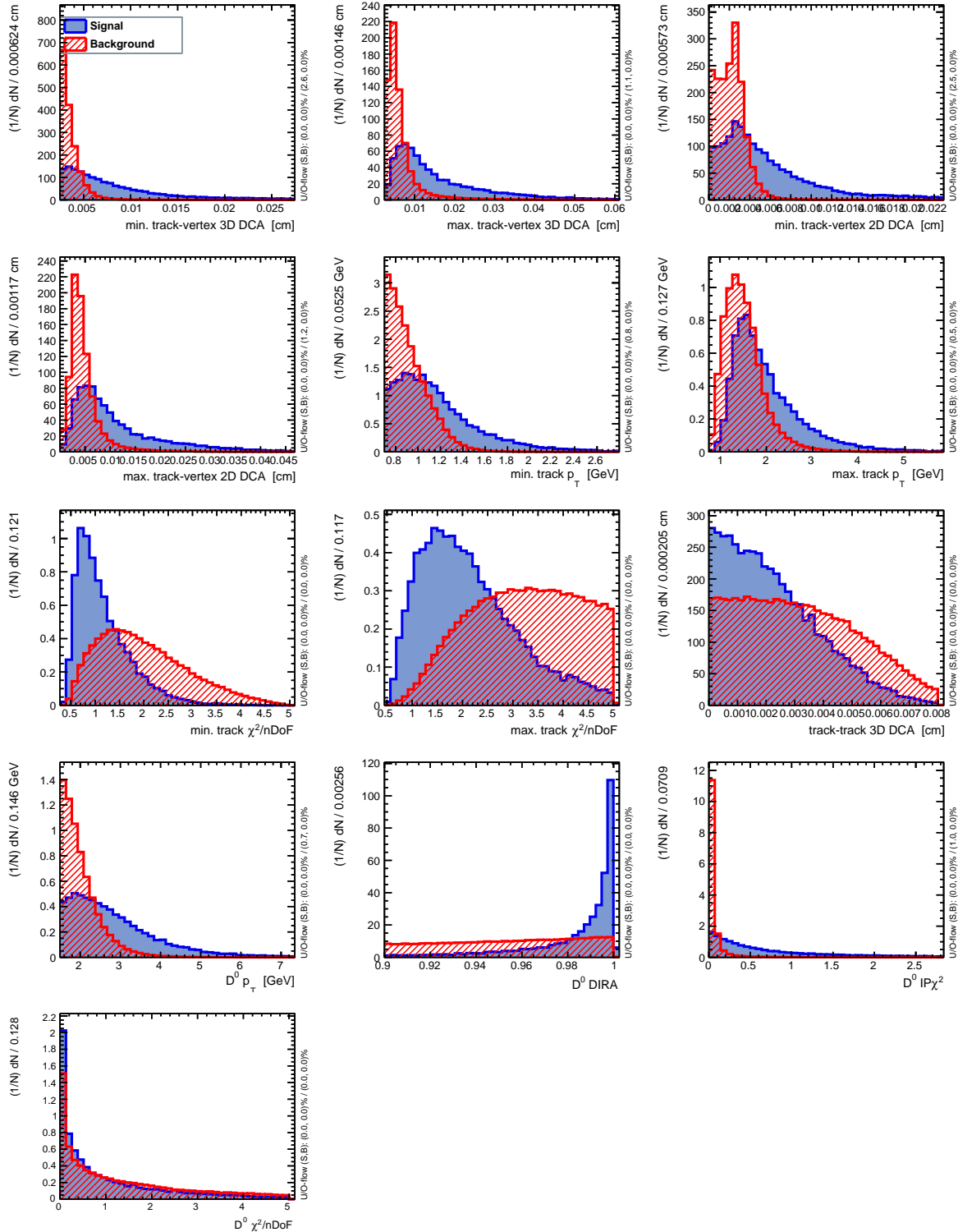


Figure 16: Input variable distributions of the signal (blue) and background (red) samples used to train the machine learning algorithms. Variables from top to bottom, left to right: minimum track IP, min. track  $p_T$ , max. track  $\chi^2$  per no. of degrees of freedom, track-track DCA,  $D^0 p_T$ ,  $D^0$  DIRA,  $D^0$  IP $\chi^2$ ,  $D^0 \chi^2$  per no. of degrees of freedom.

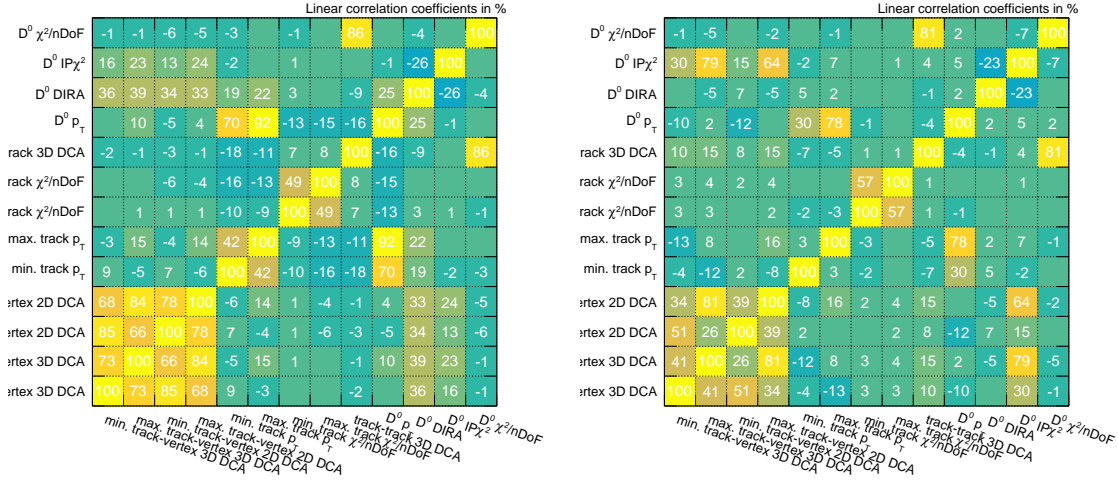


Figure 17: Input variable correlation coefficients of the signal (left) and background (right) samples used to train the machine learning algorithms.

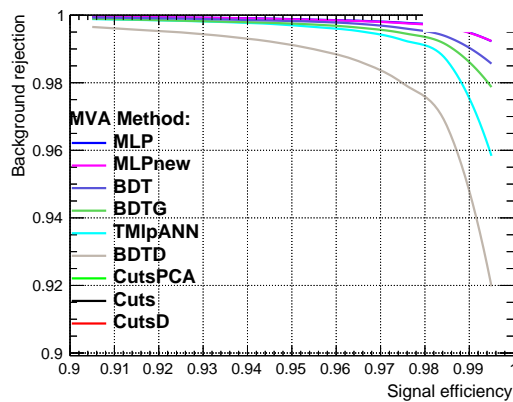


Figure 18: Receiver operating characteristic plot from the machine learning study.

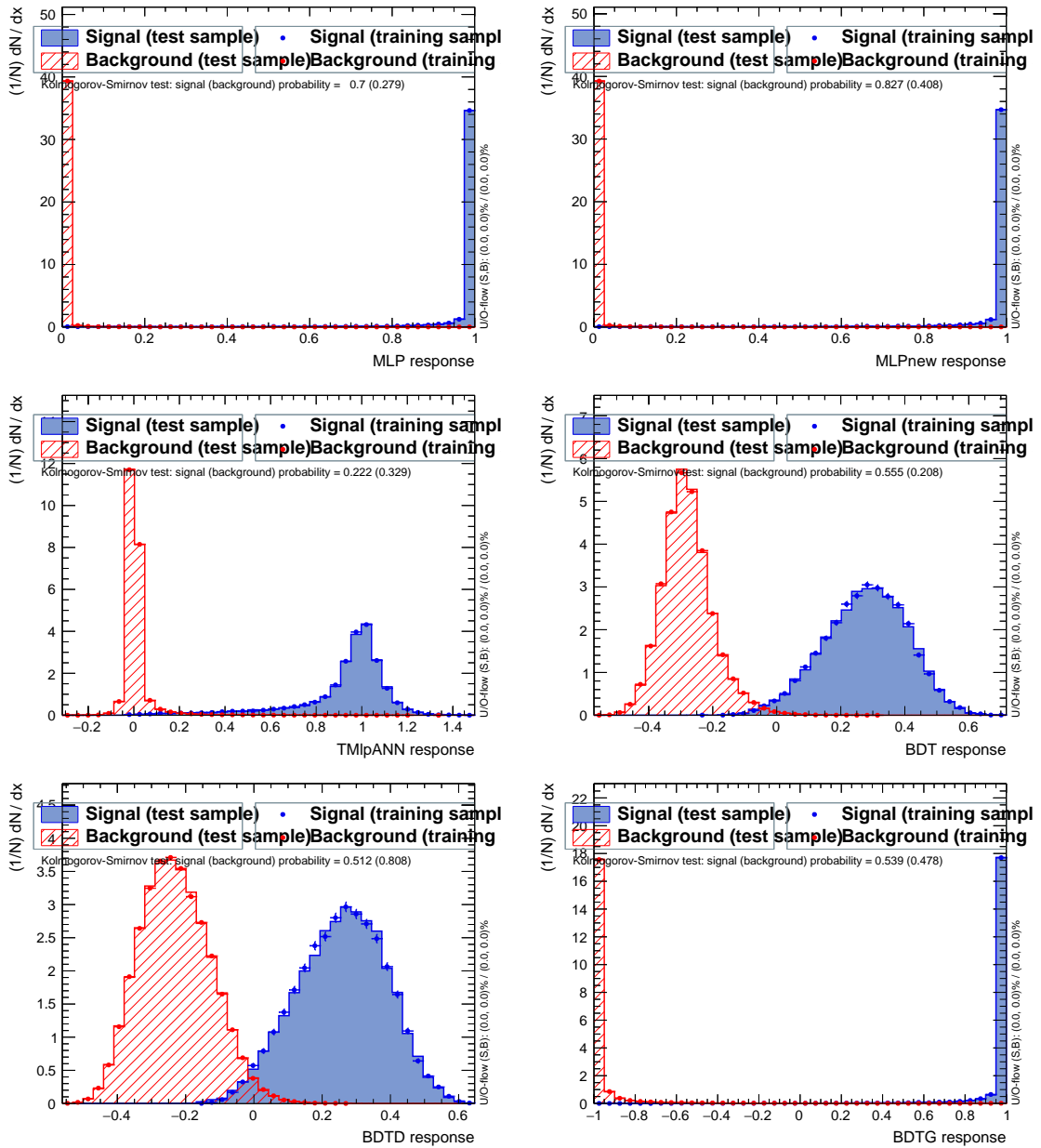


Figure 19: Classifier output distributions for the machine learning algorithms that were trained. From top to bottom, left to right: the multi-layer perceptron, a second multi-layer perceptron, ROOT's own neural net, a boosted decision tree, a boosted decision tree with decorrelation and adaptive boost, and a boosted decision tree with gradient boosting.

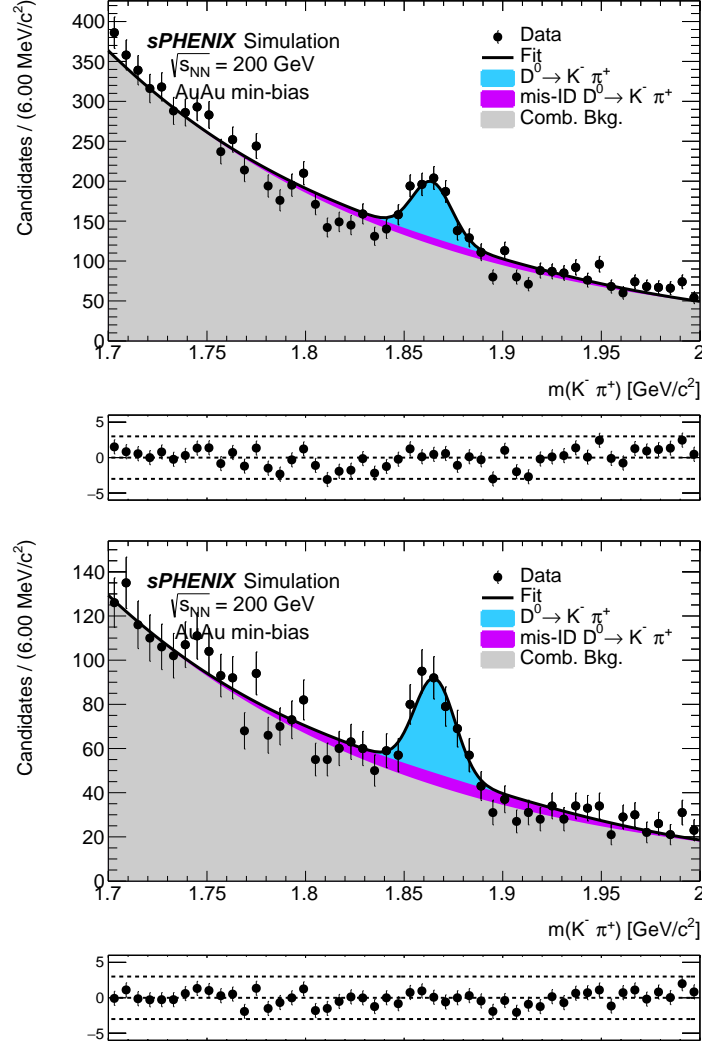


Figure 20: Fit to the  $K^- \pi^+$  invariant mass distribution using simulated Au+Au events with the baseline selection. Machine learning algorithms were trained to reject more background and this fit was performed using a multi-layer perceptron with an algorithm response  $\geq 0.99$  (top) and  $\geq 0.999$  (bottom).

## 514 C *s*Weighting

515 It is often the case that data distributions are composed of a composite of sources such a  $D^0$   
 516 signal and combinatorial background. Often we want to look at the distribution from a single  
 517 source and removing the other sources can be difficult. One method for doing this is to cut on  
 518 variables that have a large separation between sources but this can lead to contamination and  
 519 biases in the final distribution due to overlaps. Another method, known as *sWeighting* [25],  
 520 involves maximising a likelihood function for a discriminating variable (such as a candidate's  
 521 mass) where each candidate's contribution to a specific class is determined by calculating  
 522 the likelihood with and without that event then assigning an event probability where the  
 523 probabilities of an event belonging to that class are required to sum to one over all classes  
 524 in the model. These weights can then be applied to a control variable (such as a particle's  
 525  $p_T$ ) assuming there is no correlation between the two variables.

526 To investigate whether *sWeights* can be used to improve the low  $p_T$  reach of the sample,  
 527 the fit was redone using the baseline selection with an ML response  $\geq 0.999$  then the weights  
 528 were calculated for each candidate to be a signal or background event. For simplicity, the  
 529 second Gaussian that models the mis-ID was removed. The weight for each candidate to come  
 530 from a  $D^0 \rightarrow K^- \pi^+$  decay was then added to the histogram of the mother's  $p_T$  distribution  
 531 and compared to the unweighted sample. Of 3149 candidates in the total fit, the yield of  
 532 correctly reconstructed  $D^0$  was measured to be  $209 \pm 26$ . The fit to the invariant mass  
 533 distribution is given in Figure 21 and the comparison of the weighted and unweighted  $p_T$   
 534 distributions is given in Figure 22. From the weighted sample, it appears that we can reach  
 535 a  $p_T$  as low as 1.5 GeV which is the baseline cut.

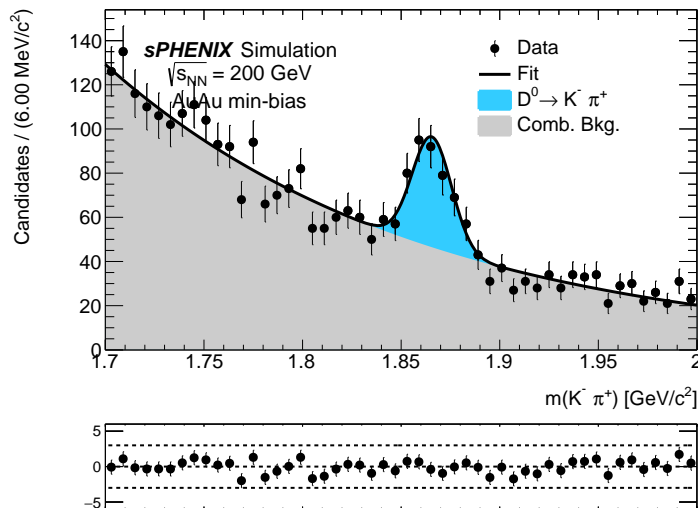


Figure 21: Fit to the  $K^- \pi^+$  invariant mass distribution using simulated Au+Au events with the baseline selection and an MLP response  $\geq 0.999$ . The mis-ID'd  $D^0$  model was removed in this fit and the fitter was required to *sWeight* each candidate.

536 ROOFIT has an internal class that is capable of calculating the *sWeights* so it is a rea-  
 537 sonably simple addition to add these weights to an nTuple. It should be noted that there  
 538 are some requirements:



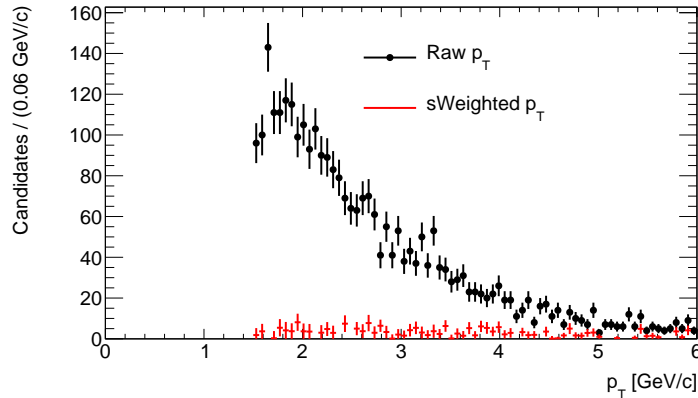


Figure 22:  $p_T$  distribution of the  $D^0$  candidates with *sWeighting* in red and without *sWeighting* in black.

- 539 1. You must have a fit value for each component of your fit. This means you can't have a  
 540 fit fraction for the signal and then subtract that from 1 to get the background fraction.
- 541 2. You must fit for the yield of each component of the fit, not the fraction. This is because  
 542 the log-likelihood minimisation requires a value that is not less than 1.

543 An example code follows,

```
stringstream cutStream;
cutStream << "1.7 <= DO_mass && DO_mass <= 2.0";
TCut masscut = cutStream.str().c_str();
TFile* dataFile = new TFile("myInputFile.root");
TTree* dataTree = (TTree*)dataFile->Get("DecayTree");

string datasWeight = inputFile.substr(0, inputFile.size()-5) + "_sWeighted.root";
TFile* sWeightedDataFile = new TFile(datasWeight.c_str(), "RECREATE");
TTree* dataSWTree = dataTree->CopyTree(masscut);
TTree* sWeightedDataTree = dataSWTree->CloneTree(-1);

RooRealVar mass(branch.c_str(), "mass", minMass, maxMass);
RooDataSet dataSet(branch.c_str(), "data", mass, Import(*sWeightedDataTree));

/*
 * Signal Model
 */
RooRealVar mean("mean", "mean", 1.865, 1.835, 1.875);
RooRealVar sigma("sigma", "sigma", 0.006, 1e-3, 0.030);
RooGaussian DO("DO", "DO", mass, mean, sigma);

RooRealVar fSig("fSig", "fSig", 0.1*dataSet.numEntries(), 0, 2*dataSet.numEntries());
RooRealVar fBkg("fBkg", "fBkg", 0.9*dataSet.numEntries(), 0, 2*dataSet.numEntries());
```

```

/*
 * Background Model
 */
RooRealVar expConst("expConst", "expConst", -10, -1e2, 0.);
RooExponential background("background", "background", mass, expConst);

/*
 * Fitting to the data
 */
RooArgList fitModellist(D0, background), fitFracList(fSig, fBkg);

RooAddPdf model("model", "model", fitModellist, fitFracList);
model.fitTo(dataSet);

RooStats::SPlot* sData = new RooStats::SPlot("sData", "An sPlot", dataSet,
                                             &model, RooArgList(fSig, fBkg));
double sig_sw;
TBranch* b_sig_sw = sWeightedDataTree->Branch("sWeight", &sig_sw, "sWeight/D");

std::cout << "Check sWeights:" << std::endl;
std::cout << "Yield of signal is " << fSig.getVal()
           << ". From sWeights it is " << sData->GetYieldFromSWeight("fSig")
           << std::endl;

for (int i = 0; i < dataSet.numEntries(); ++i)
{
    if (i < 5)
    {
        std::cout << "Signal Weight = " << sData->GetSWeight(i, "fSig")
                  << ", Background Weight = " << sData->GetSWeight(i, "fBkg")
                  << ", Total Weight = " << sData->GetSumOfEventSWeight(i)
                  << std::endl;
    }

    const RooArgSet* row = dataSet.get(i);
    sig_sw = (double) row->getRealValue("fSig_sw");
    b_sig_sw->Fill();
}

sWeightedDataFile->Write();
sWeightedDataFile->Close();

```