

Library Event Matching event classification algorithm for electron neutrino interactions in the NOvA detectors

C. Backhouse, R. B. Patterson

Lauritsen Laboratory, California Institute of Technology, Pasadena, California 91125, USA

Abstract

We describe the Library Event Matching classification algorithm implemented for use in the NOvA $\nu_\mu \rightarrow \nu_e$ oscillation measurement. Library Event Matching, developed in a different form by the earlier MINOS experiment, is a powerful approach in which input trial events are compared to a large library of simulated events to find those that best match the input event. A key feature of the algorithm is that the comparisons are based on all the information available in the event, as opposed to higher-level derived quantities. The final event classifier is formed by examining the details of the best-matched library events. We discuss the concept, definition, optimization, and broader applications of the algorithm as implemented here. Library Event Matching is well-suited to the monolithic, segmented detectors of NOvA and thus provides a powerful technique for event discrimination.

Keywords: library matching, classification algorithm, particle identification, NOvA

1. Introduction

Classifying images into a small number of categories is a common task in scientific and industrial fields. In particle physics, this task usually involves interpreting particle detector data to determine the type of particles, interactions, or decays present. Given the sheer volume of information that can be collected, the data is often first reduced to a set of derived quantities by running algorithms that pull out key features: clusters, tracks, showers, jets, etc. While this form of lossy compression is acceptable in some applications, it is worth exploring whether a classification scheme that uses all of the available information is feasible, even in cases where the data volume is high.

In this article we describe such a classification scheme developed to categorize neutrino scattering events recorded in the NOvA detectors. In the Library Event Matching (LEM) algorithm, a trial event of unknown type is compared to a large number of known “library” events to find those events that are most similar to the trial event. The properties of those best-matched library events reveal the likely nature of the trial event. A distinguishing feature of LEM is that the comparisons are made using the energy depositions directly, to avoid any information loss from calculating higher-level variables. This fundamental philosophy of LEM was developed within the MINOS collaboration for its own neutrino event categorization needs [1, 2, 3, 4]. The LEM version described in this article has substantial differences from its predecessor, many of which are motivated by the higher spatial resolution of the NOvA detectors.

While we use NOvA as our case study, the approach discussed is generalizable and could be usefully applied to any highly segmented detector, from hadron calorimeters determining jet multiplicity to cubic kilometer arrays collecting neutri-

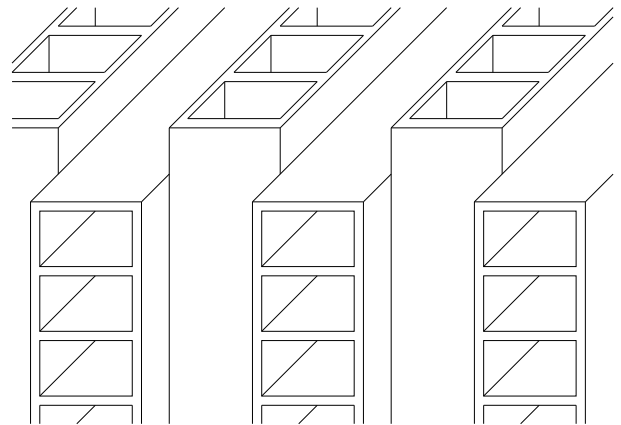


Figure 1: A sketch of the structure of the NOvA detectors. 4 cm \times 6 cm cells run the length of each 16 m \times 16 m plane. The alternating vertical and horizontal orientations can be seen. They are filled with liquid scintillator and each contains a looped wavelength-shifting fiber (not shown), as described in the text. This cut-away sketch is diagrammatic only. The real cells have rounded corners and the ends of the cells are capped for instrumentation and oil containment purposes. The neutrino beam is incident from the left.

nos from astrophysical sources. As with many machine learning algorithms, LEM requires a large number of known examples from each classification category. In particle physics applications, these would typically come either from an advanced Monte Carlo simulation or from calibration sources.

2. The NOvA experiment

The NOvA (NuMI Off-axis ν_e Appearance) experiment studies the phenomenon of neutrino flavor oscillation [5]. Neutrinos

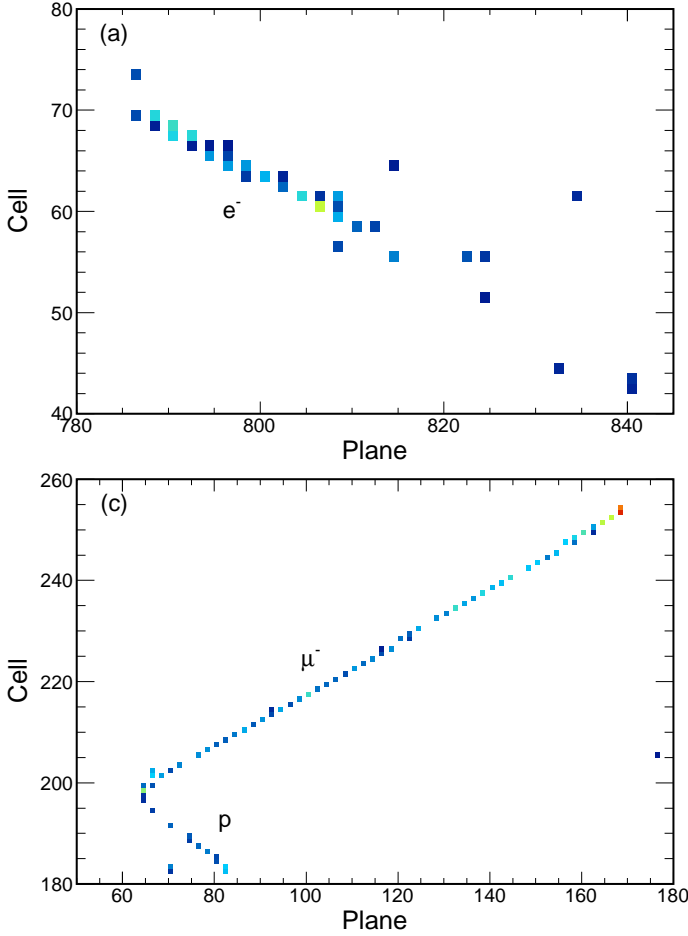


Figure 2: Example simulated events in the NOvA detectors. Only one of the two views is shown in each case. Each box represents one cell and is positioned according to its plane number (horizontal axis) and cell number (vertical axis). The color scale indicates the charge deposited in photoelectrons, and is common to all three panels. (a) A ν_e CC event, with the electron-induced electromagnetic shower clearly visible. (b) A neutral current event with a π^0 . The upper track is due to a proton. This event shows that the two showers from $\pi^0 \rightarrow \gamma\gamma$ are not always distinct. (c) A ν_μ CC event, with the usual tell-tale long, straight muon track. Note that the axis ranges are approximately doubled for this panel relative to the first two.

produced by the NuMI beamline at Fermilab [6] are observed by a Near Detector on the Fermilab site and by a Far Detector of identical construction located 810 km downstream in Ash River, Minnesota. For the purposes of this article, the neutrino oscillation mode of interest is $\nu_\mu \rightarrow \nu_e$, and the goal of the classification algorithm is to obtain a sample of electron neutrino interactions in the Far Detector with the highest possible efficiency and purity.

The NOvA detectors are constructed from long PVC cells filled with scintillator-doped mineral oil. Each of the Far Detector’s 344,064 cells is 16 m long with rectangular cross section 4 cm \times 6 cm. A loop of wavelength-shifting fiber runs the length of each cell, with both ends of the fiber terminating at one pixel of a 32-pixel APD array. The body of the 14-kiloton detector consists of 896 layers, or “planes”, each with 384 cells. Each plane is 16 m \times 16 m square, and the depth of the detector along the beam direction is 60 m. Alternate planes are aligned vertically and horizontally so that three-dimensional information can be obtained through combination of the two “views”. The detector has unprecedented granularity for its size, with one radiation length (38 cm) extending over many cells, to give a detailed view of neutrino-induced electromagnetic showers. Figure 1 shows a cut-away diagram of the detector’s construction.

The signal for the $\nu_\mu \rightarrow \nu_e$ oscillation analysis in NOvA is ν_e charged-current (CC) scattering, which yields a high-energy electron in the final state that allows one to tag the incident

neutrino’s flavor. In the 1 to 3 GeV energy range of NOvA, this electron will be accompanied, with similar probabilities, by a proton (quasi-elastic scattering), a nucleon plus a pion (resonant scattering), or a richer hadronic shower (deep inelastic scattering). While nuclear effects blur these crisp definitions, these three scattering types are useful for conveying the variety of shapes that signal events in NOvA can take. The ~ 1 GeV electron in the final state produces an electromagnetic shower in the detector that has a width of a few cells and runs longitudinally an average distance of 2.5 m (40 planes). Figure 2a shows a simulated ν_e CC interaction in the NOvA Far Detector.

The primary mis-identification background comes from neutral-current (NC) interactions, particularly those where the recoil hadronic system contains a π^0 . The π^0 decays quickly to two photons, each of which induces an electromagnetic shower that is essentially indistinguishable from an electron-induced shower. NC π^0 events, taken as a whole, look sufficiently different from signal ν_e CC events that we can reject them well, but the differences are sometimes obscured:

- The presence of two electromagnetic showers, rather than one, can reveal a π^0 in the final state. However, if one of the showers has low energy or overlaps the other in the detector, it can be missed.
- Photon-induced showers are separated from the neutrino interaction point due to the distance traveled by the photon prior to its conversion. This gap is a tell-tale sign of

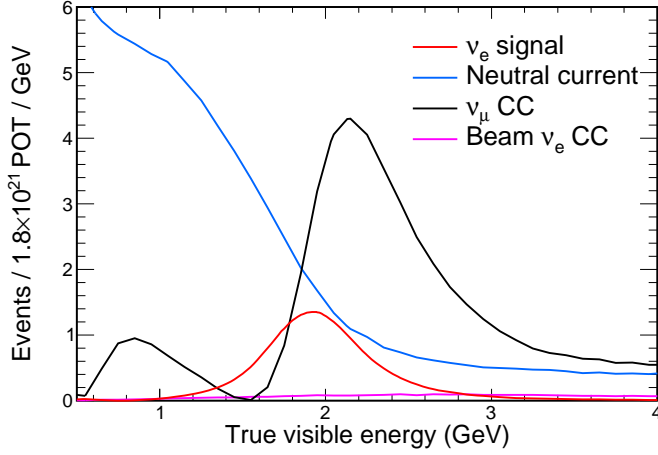


Figure 3: Signal and background distributions of visible energy expected in the Far Detector sample. The effect of neutrino oscillations is included. Visible energy is defined as the incident neutrino energy except in the case of neutral current events where the outgoing neutrino energy is subtracted. The ν_e signal to be identified by LEM is shown in red. The neutral current, ν_μ charged current, and intrinsic beam ν_e charged current components are blue, black, and magenta respectively.

a photon, but in some cases the gap will be too small to resolve. The conversion length in NOvA is 50 cm.

- Photon-induced showers begin with two particles (an electron/positron pair) rather than one, but these cases can end up indistinguishable given the energy resolution of the detector.
- The energy lost to the outgoing neutrino in NC scattering leads to reconstructed energies lower than those of signal events. However, interactions from a sufficiently high-energy neutrino or with a large energy transfer can fall in the signal region of 1 to 3 GeV reconstructed energy.

Figure 2b shows a simulated NC event with a π^0 .

Additional background comes from ν_μ CC scattering, which produces a muon in the final state. The muon leaves a long track of activity in the detector with a characteristic energy deposition per unit pathlength. These are readily removed from the sample due to the clear muon track except in cases where the muon is low in energy or is lost amongst other activity. In these cases, the background is similar to NC interactions, with neutral pions playing the same role. Figure 2c shows a ν_μ CC example.

The NuMI beam also includes a 2% contamination of ν_e . These ν_e interact identically to the ν_e from oscillations and thus constitute a background to the $\nu_\mu \rightarrow \nu_e$ oscillation measurement. However, their rate is low and their energies are somewhat higher. Figure 3 illustrates the energy differences among all the event classes before any selection cuts have been applied.

Since the ν_e CC signal falls within a known energy range, we can safely remove lower and higher energy events up front. For all figures and tables that follow, we require events to have reconstructed visible energies between 0.5 GeV and 4 GeV.

3. Library Event Matching concept

At the heart of the LEM algorithm is the comparison of each unknown trial event to a large number of known library events, with the comparisons based on low-level information collected by the detector. For NOvA, this means using the calibrated energy depositions in all the detector cells directly rather than forming higher-level objects such as showers and tracks from those.

Once the very best matches are found (here, the best 0.0001% of all library events), their known properties are used to estimate the properties of the trial event. In the simplest version of LEM, the fraction of the best matches that are signal events can be used as the discriminant. Appendix A.1 discusses the relationship between LEM and other machine learning techniques.

3.1. The matching metric: motivation

When comparing two events, a metric is needed to quantify how similar they are. It is instructive to look at the MINOS case briefly, as the situation there is somewhat simpler [1, 2, 3, 4].

The MINOS detector has a segmented structure analogous to that of the NOvA detector, but the effective spatial resolution for events of interest is significantly lower. A ν_e CC signal event in MINOS involves only a couple dozen active “strips” (the analogue of NOvA’s cells), and these active strips are clustered in a relatively compact pattern. Thus, two events with the same underlying particle kinematics have a good chance of having identical (or near-identical) arrangements of active strips. The read-out electronics report the number of photoelectrons detected in each active strip. Since this charge measurement suffers from shot noise (typical charge: ~ 8 photoelectrons), strips with identical energy depositions may report different charges. The level of difference is governed by Poisson statistics.

These details guided the form of matching metric used by MINOS, which can be thought of as the likelihood \mathcal{L} that the two events’ recorded charges represent the same underlying energy depositions:

$$\log \mathcal{L} = \sum_i^{\text{strips}} \log \left[\int P(a_i|\lambda)P(b_i|\lambda)d\lambda \right], \quad (1)$$

where a_i is number of photoelectrons registered by the i^{th} strip of event A, b_i is the same for event B, $P(n|\lambda)$ is the Poisson probability of observing n given mean λ , and the sum runs over all strips active in at least one of the events. A higher $\log \mathcal{L}$ for a pair of events means a better match. Before \mathcal{L} is calculated, the events, which in general occur in different parts of the detector, are spatially aligned by shifting them so that their charge-weighted mean strip positions, rounded to the nearest strip, overlap.

In the MINOS metric \mathcal{L} , displaced energy depositions in the two events do not get their charges directly compared. To obtain good matches for a trial event, the library must be large enough to span minor variations in active strip positions for nominally equivalent events. This is possible in MINOS given the limited spatial resolution of the detectors for ν_e CC events. That is, the library can be expected to give reasonable coverage of all

possibilities. Requiring exact *charge* agreement across the ~ 20 active strips, though, would be combinatorically overwhelming. The Poisson factors take care of this, with acceptably different charges able to contribute appropriately to the match score.

The NOvA detectors are significantly more finely-grained than those of MINOS. This makes event discrimination easier in principle since more details are visible, but it makes the above matching metric impractical. It is much less likely that “equivalent” activity in the trial and library events will fall on the same cells. What is needed is a matching metric that rewards activity in nearby cells without requiring them to lie directly on top of one another. A library event identical to the trial event should still be a perfect match, but events with similar charges offset by a cell or so should still score well.

The metric we use draws its motivation from electrostatics. Two Coulomb charge distributions of similar shape, but with opposite signs, will have a low electrostatic potential energy when overlaid and examined together, as the attraction between the opposite signed charges counters the internal repulsion of the like-signed charges. Two overlaid charge distributions with dissimilar shape suffer the internal repulsion but lack the benefit of mutual attraction, leading to a large potential energy. Given the electrostatic analogue to what follows, we use “energy” to refer to the LEM match score for the remainder of the article unless otherwise stated. Lower energies correspond to better matches.

The match energy is defined as

$$E = E_A + E_B + E_{AB}, \quad (2)$$

where E_A is the self-energy (repulsion) of event A’s charges, E_B is the self-energy of event B’s charges, and E_{AB} is the (negative) energy due to the A/B attraction. The charges are taken to be the recorded energy depositions in the NOvA cells. Treating the electrostatic analogue as exact for a moment, the self-energy terms are given by

$$E_A = \frac{1}{2} \sum_{ij}^{\text{cells}} \frac{a_i a_j}{r_{ij}}, \quad E_B = \frac{1}{2} \sum_{ij}^{\text{cells}} \frac{b_i b_j}{r_{ij}}, \quad (3)$$

with a_i (b_i) the recorded deposition in the i^{th} cell of event A (event B) and with r_{ij} the distance between cells i and j . The $r_{ij} = 0$ case is handled again with an electrostatic analogue by distributing all charges uniformly across their individual cells. (See Appendix A.2.)

The interaction term is given by

$$E_{AB} = - \sum_{ij}^{\text{cells}} \frac{a_i b_j}{r_{ij}}. \quad (4)$$

Before evaluating this sum, the events are globally aligned with one another according to a separately reconstructed interaction vertex.¹

A perfect match, in which events A and B have identical depositions in identical cell positions, would yield $E=0$. A poorly

matched pair with charges far away from one another will have large energy $E \approx E_A + E_B$.

Eq. (4) can be recast in terms of one set of charges embedded in the field of the other:

$$E_{AB} = - \sum_i^{\text{cells}} a_i V_i \quad (5)$$

$$\text{where } V_i = \sum_j^{\text{cells}} \frac{b_j}{r_{ij}}. \quad (6)$$

The advantage of this formulation is that V can be precalculated for each trial event, along with the self-energies of the trial and library events. When matching against a large number of library events using (5), the complexity is linear in the number of charges rather than requiring a double sum over both trial and library charges.

3.2. The matching metric in NOvA

While the NOvA matching metric is inspired by electrostatics, there is no reason to expect that the precise form above will yield the best sensitivity. We incorporate the following generalizations.

- Above, r_{ij} is calculated as the Euclidean distance in terms of the number of planes Δp_{ij} and number of cells Δc_{ij} . However, NOvA events are boosted forward and cover many planes longitudinally but relatively few cells transversely, so we assign different relative importance to separations in the two directions.
- The r^{-1} falloff with distance is generalized to $r^{-\alpha}$.
- The importance of larger charges relative to smaller ones is adjusted by raising all charges to a power β .

The resulting form of the matching metric still follows Eq. (2), but the self-energy and interaction terms are now given by

$$E_A = \frac{1}{2} \sum_{ij}^{\text{cells}} a_i^\beta T_{ij} a_j^\beta, \quad E_B = \frac{1}{2} \sum_{ij}^{\text{cells}} b_i^\beta T_{ij} b_j^\beta \quad (7)$$

$$E_{AB} = - \sum_i^{\text{cells}} a_i^\beta U_i \quad (8)$$

with the transfer matrix T_{ij} and field U_i given by

$$T_{ij} = \left(\frac{\Delta p_{ij}^2}{\sigma_p^2} + \frac{\Delta c_{ij}^2}{\sigma_c^2} \right)^{-\alpha/2} \quad (9)$$

$$U_i = \sum_j^{\text{cells}} T_{ij} b_j^\beta. \quad (10)$$

The electrostatics version is recovered by setting

$$\sigma_p = \sigma_c = \alpha = \beta = 1. \quad (11)$$

¹Alignment by charge-weighted mean cell position was also studied and gives similar classification performance.

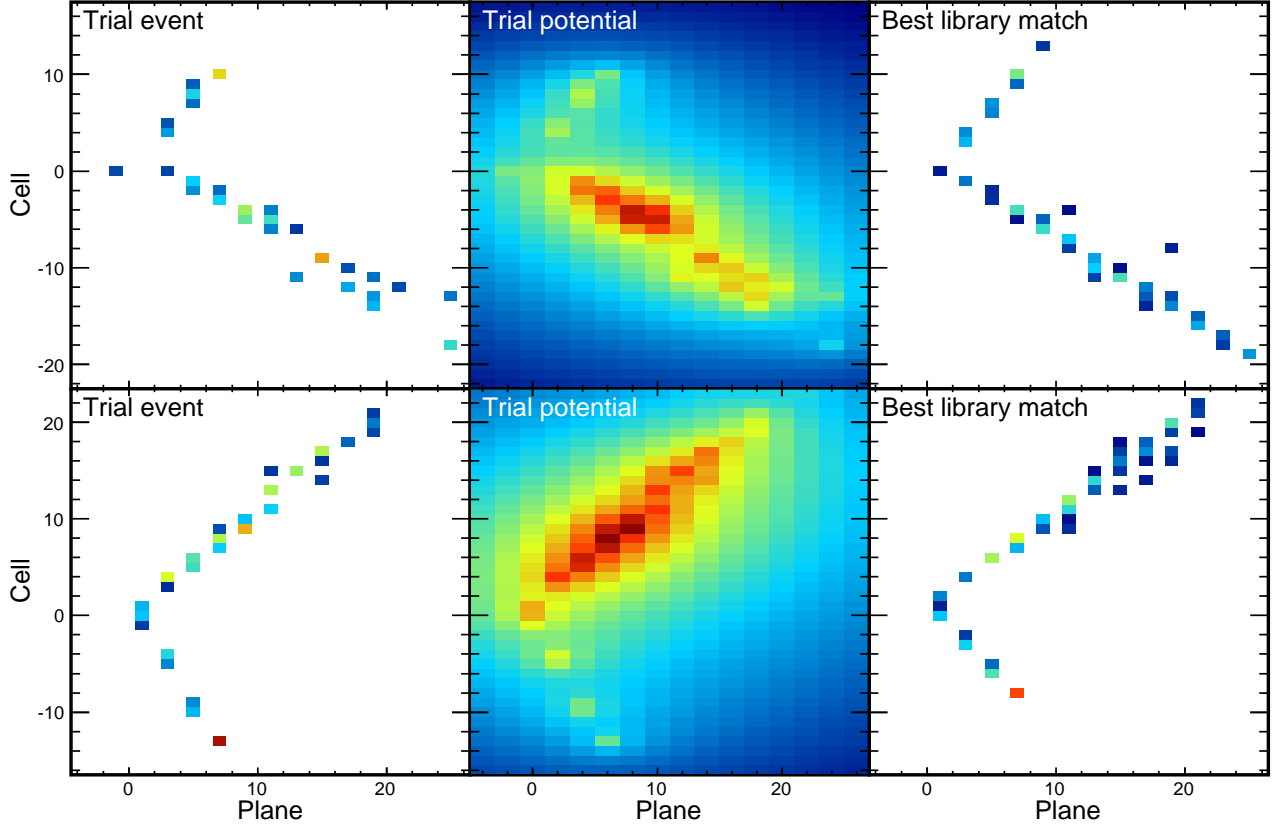


Figure 4: Example of LEM matching. On the left is a trial ν_e CC event, on the right the best match found. The central panels shows the potential U in which the library events are placed in order to calculate the match energy. The upper panels show one view, and the lower panels show the other.

We ran toy experiments with different values of these parameters and calculated a figure-of-merit for each to optimize performance. The parameters chosen were:

$$\sigma_p = 0.286 \quad (12)$$

$$\sigma_c = 0.095 \quad (13)$$

$$\alpha = 0.25 \quad (14)$$

$$\beta = 0.5 \quad (15)$$

The first two parameters validate the intuition that transverse differences should be considered more significant than longitudinal ones. The third parameter specifies a $1/\sqrt{r}$ falloff with distance, slower than the electrostatic analogue. For β , note that the simple presence or absence of activity in a cell conveys information regardless of its charge. Having $0 < \beta < 1$ moves the metric towards this binary “on/off” interpretation and away from a charge-proportional weighting.

4. The library

The library consists of 77M simulated neutrino events, of which 18M are signal ν_e CC events, 29M are background ν_μ CC and NC events, and 30M are π^0 -enriched NC background events. Each trial event that LEM classifies is compared to these 77M events to find the 1,000 library events that are most sim-

ilar to it, as quantified by the metric above.² Figure 4 shows an example trial event along with its event potential U and its best-matched library event.

The library events are generated ahead of time using the full NOvA Monte Carlo simulation chain including realistic neutrino flux, cross sections, and detector components. The flux is calculated using a FLUKA/FLUGG implementation of the beamline elements [7], the neutrino interactions are simulated by GENIE [8], and particle propagation through the detector geometry is handled by GEANT4 [9]. Simulated energy depositions in the liquid scintillator are converted into expected signals by NOvA electronics and data acquisition simulation code. The registered signals are corrected for light attenuation in the cells’ fibers using standard NOvA calibration procedures.

NC events containing neutral pions are the dominant mis-identification background owing to the electromagnetic showers from $\pi^0 \rightarrow \gamma\gamma$. Thus, we supplement the base background library sample with a π^0 -enriched library sample. To build this enriched sample, we apply a cut that selects out only those neutral current events with a π^0 present in the final state as reported by GENIE.

The library events are generated according to the expected ν_μ flux (for background) or a 100% $\nu_\mu \rightarrow \nu_e$ transmutation (for signal), without regard to any actual probabilities for neutrino

²This statement is modified in Sec. 7.1 when we discuss speed optimizations.

flavor change. Oscillations are introduced into the library later by event weighting. This is discussed in Sec. 5 below. Appendix A.3 describes the oscillation probabilities used.

While increasing the library size beyond the 77M events would provide incremental improvement in classification performance, we observe that these gains enter logarithmically with the number of library events once the library is sufficiently large. In an earlier version of the algorithm, we found that doubling the library size provided only 1% gain in physics sensitivity. In light of the computational requirements discussed in Section 7, additional library events are not worthwhile for our application.

4.1. Event flipping

To good approximation, flipping an event transversely in one or both views produces an equally valid event. We use such flipping to effectively quadruple the size of the library when the matching is performed. Each library event is used in each of the four possible configurations, and the best of the four is retained. This symmetry is not quite perfect in the NOvA detectors. Attenuation in the readout fibers leads to subtly different charge resolutions and threshold effects on transversely opposing sides of an event, and NuMI neutrinos at the Far Detector enter at a 3° upwards angle. Nevertheless, the best-scoring matches come from the four possible flipped configurations with nearly equal probability: 26% from unflipped events, 50% from events with either one of the two views flipped, and 24% from events with both views flipped.

5. Decision tree

As library size increases, the fraction of an event's best matches that are truly signal tends toward the probability that the trial event itself is signal. Further, all of the information available in the trial event is used when determining this probability. It is in this sense that LEM is optimal.

For a library of finite and practical size, though, this signal fraction alone does not contain the full information extractable. Other statistics constructed from the details of the best matches may, for example, indicate that the matches are drawn from an area of sparse library coverage and are thus less reliable. The most powerful approach given a finite library is to construct several statistics describing the matches and to feed these into one of the standard multivariate analysis techniques to extract the final classifier. In LEM, five variables are constructed from the 1,000 best library matches and are used as inputs to a decision tree, along with the calorimetric energy of the trial event as a sixth input.

5.1. Weighted fraction of signal matches

The basic quantity measuring what fraction of the best matches are signal events can be improved upon by weighting up the truly best matches over the lesser ones when calculating the signal fraction. We use the weighting

$$w'_n = \exp\left(-\lambda \left(\frac{E_n}{E_{1000}}\right)^\gamma\right), \quad (16)$$

where n is the match index, E_n is the energy of the n^{th} best match for the trial event, and E_{1000} is the energy of the final (1000th) best match. The optimized values used for λ and γ in NOvA are

$$\lambda = 6.67 \quad (17)$$

$$\gamma = 10. \quad (18)$$

The typical ratio of weights w'_{1000}/w'_1 is $\sim 0.1\%$, indicating that the most important matches are captured within the first thousand.

In practice, the weight must also include the oscillation probabilities alluded to earlier:

$$w_n = w'_n P_n^{\text{osc}}, \quad (19)$$

where P_n^{osc} is the oscillation probability of match n , as described in Appendix A.3.

All sums below that are indexed by n run over the match list. For notational convenience we also define $W \equiv \sum_n w_n$. This weighting scheme is used for all five quantities formed from the best-match list. The first is the weighted fraction of signal matches,

$$f_{\text{sig}} = \frac{1}{W} \sum_{n, \text{sig}} w_n, \quad (20)$$

where this sum includes only those terms due to signal matches.

5.2. Mean hadronic y

Signal events in which the outgoing electron carries only a small fraction of the incident neutrino's energy will look very much like NC background events. The kinematic quantity y (or rather, $1-y$) measures this fraction: $1-y = K_e/K_\nu$, where we've used K_e and K_ν as the outgoing and incoming lepton energies to avoid confusion with the match energies E . If a trial event matches well to signal events with high y , this can suggest that the trial event is in fact a high- y NC event. A second input is the mean y for the best matches:

$$\langle y \rangle = \frac{1}{W} \sum_n w_n y_n. \quad (21)$$

5.3. Mean matched charge fraction

Matched charge fraction is an independent measure of the quality of the library matches, separate from the match energy. For each trial/match pair, this is the quantity of charge that has a counterpart on identical cells in the two events divided by the total charge in the two events:

$$f_Q = \frac{2 \sum_i^{\text{cells}} \min(a_i, b_i)}{\sum_i^{\text{cells}} (a_i + b_i)}. \quad (22)$$

The weighted average of the matched charge fraction over all the matches yields the next input:

$$\langle f_Q \rangle = \frac{1}{W} \sum_n w_n f_{Q,n}. \quad (23)$$

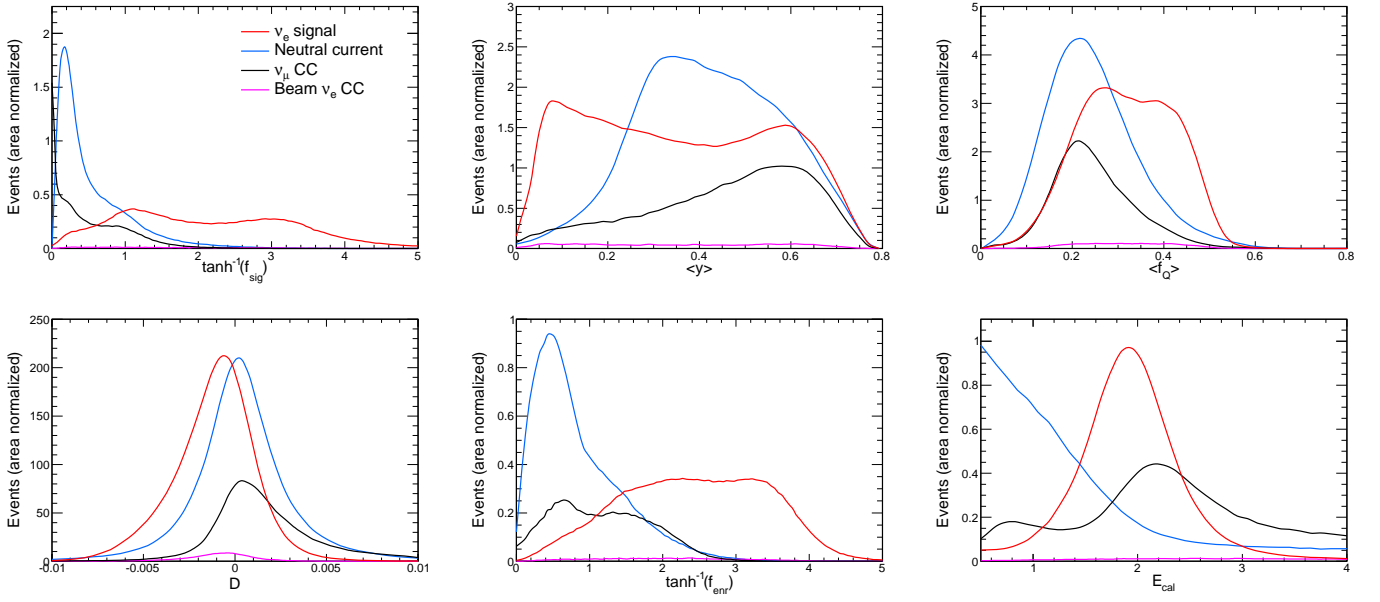


Figure 5: The six decision tree inputs described in the text. The red curves show the distribution of signal events. The blue, black, and magenta curves show the distributions of neutral current, ν_μ CC, and intrinsic ν_e CC backgrounds respectively. The ν_e signal and neutral current background are normalized to equal area. The other backgrounds are to the same scale as the neutral current curve. The signal distributions for f_{sig} and f_{enr} are very sharply peaked at 1, so we have plotted these quantities as $\tanh^{-1}(f_{\text{sig}})$ and $\tanh^{-1}(f_{\text{enr}})$ to keep the signal and background curves visible on the same vertical scale.

5.4. Match energy difference

This quantity measures whether the signal or background matches are the better matches on average. It is the difference of the weighted mean energy of each class of matches:

$$D = \frac{\sum_{n, \text{sig}} w_n E_n}{\sum_{n, \text{sig}} w_n} - \frac{\sum_{n, \text{bkg}} w_n E_n}{\sum_{n, \text{bkg}} w_n} \quad (24)$$

5.5. Enriched fraction

The final match list quantity, similar in construction to f_{sig} , is the weighted fraction of signal matches present among the signal and π^0 -enriched matches (*i.e.*, excluding the non-enriched background),

$$f_{\text{enr}} = \frac{\sum_{n, \text{sig}} w_n}{\sum_{n, \text{enr}} w_n + \sum_{n, \text{sig}} w_n} \quad (25)$$

5.6. Total calorimetric energy

NC backgrounds skew heavily to low visible energy thanks to the energy removed by the exiting neutrino. The sum of all depositions $\{a_i\}$ recorded in the trial event,

$$E_{\text{cal}} = \sum_i^{\text{cells}} a_i, \quad (26)$$

is included as a final input so that the classifier knows the prior expectations of signal and background.

5.7. Choice of a decision tree, and figure of merit

There are many multivariate techniques capable of combining these six input quantities into a single classifier output. We investigated artificial neural networks, support vector machines, and decisions trees. An ensemble decision tree yielded the best performance of the approaches tried. One problem with other techniques is that the figure of merit (f.o.m.) that, for example, artificial neural network training aims to minimize is the mean-squared-error of the classifier variable c :

$$\text{f.o.m.} = \sum_i^{\text{sig}} (1 - c)^2 + \sum_i^{\text{bkg}} c^2, \quad (27)$$

where the sums run over the signal and background training samples. However, the figure of merit relevant to an experiment measuring the magnitude of a signal excess s over a background b with Poisson fluctuations is

$$\text{f.o.m.} = \frac{s}{\sqrt{s + b}}. \quad (28)$$

If events are binned according to, say, the classifier output, the generalization is simply to sum in quadrature the significances in the individual bins:

$$\text{f.o.m.} = \sqrt{\sum_i^{\text{bins}} \frac{s_i^2}{s_i + b_i}}. \quad (29)$$

While training a decision tree classifier, if the sample is divided at each step into subsamples 1 and 2 so as to maximize

$$\frac{s_1^2}{s_1 + b_1} + \frac{s_2^2}{s_2 + b_2}, \quad (30)$$

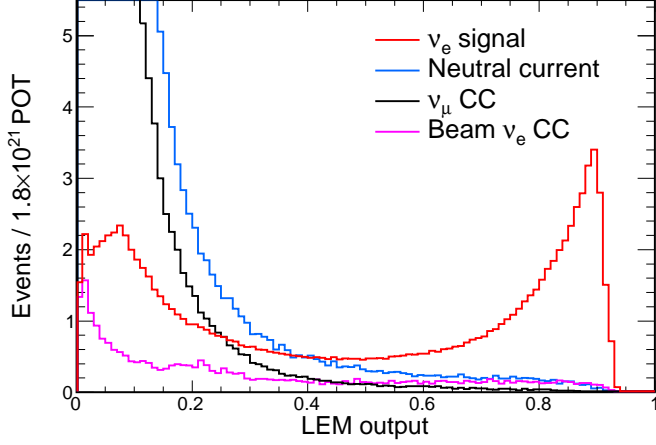


Figure 6: The distribution of the LEM output variable for ν_e CC signal events (red) compared to the background components: neutral current (blue), ν_μ CC (black) and intrinsic beam ν_e CC (magenta). In order to make the details in the signal-like region visible, the y-axis truncates much of the background peak. 95% of neutral current events and 98% of ν_μ charged current events have $\text{LEM} < 0.15$. The distributions are scaled to a nominal 3-year NuMI exposure [5] of 1.8×10^{20} protons-on-target.

then the performance of the full classifier is trivially optimized with respect to the figure of merit in Eq. (29).

The final classifier output is a voting ensemble of 1,000 decision trees each trained on a randomly chosen half of the full training sample. The ensemble technique protects against over-training, a feature that we confirmed by evaluating the classifier performance on independent control samples.

6. Classification performance

Figure 5 shows the distribution of the six input variables for all event classes in the NOvA $\nu_\mu \rightarrow \nu_e$ analysis. Figure 6 shows the final LEM classifier output. Figure 7 shows the signal efficiency and purity obtained with various cuts on the LEM output. All curves come from Monte Carlo simulation of the expected NOvA data set. We choose the cut on the LEM output variable that maximizes the figure-of-merit in Eq. (28). When applying LEM in a full experimental setting, one can fit the output distribution to gain additional discrimination power.

Table 1 shows the expected number of signal and background events selected by the optimum LEM cut. The signal efficiency is 55% for a background mis-identification rate of 2.0%. The muon track of ν_μ CC events keeps their mis-identification rate particularly low. Background beam ν_e events are selected with a lower efficiency than signal ν_e events. This is possible due to the different underlying energy spectra of the two classes. As there is no absolute metric by which to judge the performance of the LEM classification algorithm described here, we note simply that the performance shown is excellent for the physics goals of NOvA [5].

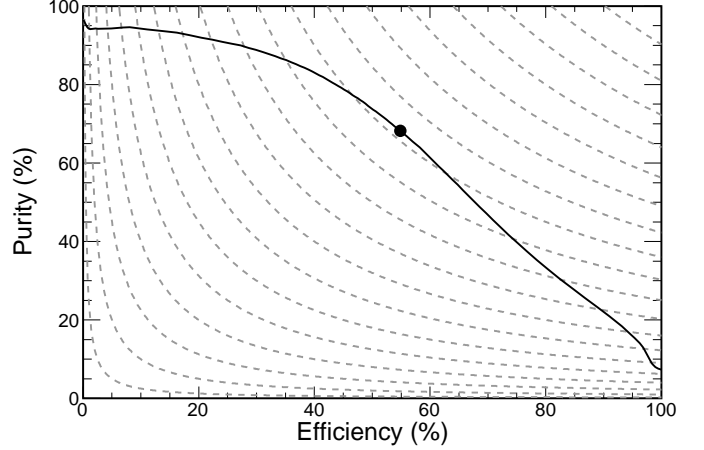


Figure 7: Efficiency and purity of the ν_e candidate sample selected by LEM for different cut positions. The dashed lines are curves of constant f.o.m. = $s / \sqrt{s + b}$, and the solid circle indicates the result of the optimum cut.

7. Computational optimization

7.1. Speed

While each individual energy calculation can be performed very quickly, classifying a single event takes some time given the large size of the library. For the NOvA application, a single event must be treated in a second or so, which is the time scale required by other steps already performed during NOvA event processing. Without specialized hardware to run the inner loop, techniques to manage the LEM matching time focus on reducing the number of energies that need to be calculated.

We achieve a significant speed-up by introducing a library “index”. If trial event A matches well to library event B , A will likely match well to other library events that are, themselves, good matches to B . Similarly, if A and B match poorly, then A will likely match poorly to library events similar to B .

A library index is formed by drawing 10,000 events uniformly from the full library and matching each of these to the full library. For each index event, a list of its 1,000,000 best-matched library events is saved. This process happens ahead of time, at library creation. When a trial event is classified, it is compared first to the 10,000 index events to find the single best-matching index event. The trial event is then compared only to the 1M sibling events of that index event, reducing the total number of energies calculated per trial event from 77,000,000 to 1,010,000 – a significant speed improvement that takes the per trial matching time from 97 s down to 1.7 s on a 2.3 GHz AMD Opteron processor. Empirically, we find that 85% of the trial event’s “true” one-thousand top matches are captured with this indexed approach, and we find no noticeable degradation in the physics performance.

7.2. Memory

The speed optimization above is what allows the use of a 77M event library. However, such a large library strains memory resources. The full library is too large (~ 53 GB each for the library and index) to read from disk for each event, yet it

	ν_e signal	Tot. bkgd.	NC	ν_μ CC	Beam ν_e CC
No selection	105	1332	734	573	25
LEM	58	27	14	4.6	7.9
Efficiency	55%	2.0%	2.0%	0.8%	32%

Table 1: Number of events expected in each event category initially and again after an optimal LEM cut assuming a nominal 3-year NuMI exposure of 1.8×10^{21} protons-on-target. The background is shown both as a total as well as broken down into NC, ν_μ CC, and intrinsic beam ν_e CC components. The bottom row shows the efficiencies for selecting events in each category. The “no selection” row and the efficiencies derived from it count only those events with reconstructed visible energy between 0.5 GeV and 4 GeV.

is larger than the typical per-core memory allocation on grid computing nodes.

Thus, the library is converted from its original high-level format into the memory representation used by a running job. This representation includes the self-energy of each event. The conversion inflates the library slightly to 131 GB, but the advantage is that it can now be shared between running processes. Each parallel matching job uses the `mmap()` system call to make the contents of this file visible in its address space. The mapping is marked read-only, so the kernel shares the pages between all the running processes. For example, on a 64-core server, the memory requirement to run 64 matching jobs is still only 131 GB, equivalent to an unshared 2 GB per core. In case of memory pressure, the kernel will discard pages, knowing that they can be retrieved from disk (that is, the library file essentially acts as swap space) although this will significantly impact performance.

8. Other information available in the match list

In addition to signal-or-background classification, the detailed truth information available in the list of best matches allows other information about the trial event to be inferred. One could extract probabilities for different interaction modes, the inelasticity, and so on, without requiring any independent reconstruction. An application that has been pursued is the estimation of the incident neutrino energy for ν_e CC events. Simply by averaging the true neutrino energies of the best signal library matches and calibrating the resulting estimator, we achieve an energy resolution of 8.8% on signal events selected by the oscillation analysis, competitive with other energy estimators in NOvA.

9. Summary

The Library Event Matching algorithm compares input trial events to a large library of known events using all the information available, making LEM an optimal classifier given a sufficiently large library. The NOvA implementation of LEM has demonstrated excellent performance in separating ν_e signal from the key backgrounds, and a few simple optimizations have maintained practical computational requirements despite the large number of library events used. Within the NOvA context, the LEM technique has potential applications from reconstruction of the hadronic system to the event energy measure

described above. More broadly, LEM can be applied to completely different particle detectors or imaging systems in an array of fields and industries, wherever one needs to classify fine-grained images of objects whose visual characteristics vary in known ways.

10. Acknowledgments

The authors thank the NOvA collaboration for use of its Monte Carlo simulation software and related tools. This work was supported by the the US DOE under award de-sc0006543.

Appendix A. Additional technical notes

A few technical notes are included in this Appendix so as not to break up the discussion in the main text.

Appendix A.1. Relation to other classification techniques

If f_{sig} and f_{enr} were calculated unweighted, then those variables would be k -nearest-neighbors classifiers, albeit with very large input vectors. With the weights w_n applied, they act as kernel density estimators. Note that

$$\sqrt{2E} = \sqrt{\sum_{ij} (a_i^\beta - b_i^\beta) T_{ij} (a_j^\beta - b_j^\beta)} \quad (\text{A.1})$$

is a metric for the space of possible event images. That is, distances defined in this way obey the triangle inequality. For a Gaussian kernel in this space one would expect $w_n \sim \exp(-E)$, which contrasts with the optimal value of $\gamma = 10$ found in practice. Similarly $\langle y \rangle$ is an estimator for the true value of y using the same kernel.

Methods exist to efficiently find nearest-neighbors in general metric spaces without having to rely on heuristics such as the library index in Section 7.1. Testing of a vantage-point tree [10] indicated its performance was affected by the curse of dimensionality. A large fraction of the nodes would have to be entered during a typical search.

Appendix A.2. Energy calculation when $r_{ij} = 0$

The transfer matrix element T_{ij} as written in Eq. (9) diverges when $i = j$ since Δp_{ii} and Δc_{ii} are zero. Thus, for nearby cell pairs ($\Delta p_{ij} \leq 5$ and $\Delta c_{ij} \leq 5$), the energy calculation is performed as if the charge is distributed uniformly over each cell, with

$$T_{ij} = \int_0^1 \int_0^1 \int_0^1 \int_0^1 [r_{ij}(x, y, u, v)]^{-\alpha} dx dy du dv, \quad (\text{A.2})$$

where (x, y) and (u, v) scan over the areas of cells i and j and where r_{ij} here is a generalization of the discrete distance used in the main text:

$$r_{ij}(x, y, u, v) = \sqrt{\left(\frac{\Delta p_{ij} + x - u}{\sigma_p}\right)^2 + \left(\frac{\Delta c_{ij} + y - v}{\sigma_c}\right)^2}. \quad (\text{A.3})$$

For more distant pairs the simplified form of the transfer matrix given in Eq. (9) is sufficient.

Appendix A.3. Neutrino oscillation weights

The retained matches are weighted according to Eq. (19), which includes the probability for flavor oscillation. The probabilities used are

$$P(\nu_\mu \rightarrow \nu_e) = \sin^2 \theta_{23} \sin^2 2\theta_{13} \sin^2 \left(\frac{1.27 \Delta m^2 L}{E} \right) \quad (\text{A.4})$$

$$P(\nu_e \rightarrow \nu_\mu) = \sin^2 \theta_{23} \sin^2 2\theta_{13} \sin^2 \left(\frac{1.27 \Delta m^2 L}{E} \right) \quad (\text{A.5})$$

$$P(\nu_\mu \rightarrow \nu_\mu) = 1 - \sin^2 2\theta_{23} \sin^2 \left(\frac{1.27 \Delta m^2 L}{E} \right) \quad (\text{A.6})$$

$$P(\nu_e \rightarrow \nu_e) = 0, \quad (\text{A.7})$$

where $L = 810$ km is the oscillation baseline, E is the neutrino energy in GeV, and the oscillation parameters are taken to be

$$\theta_{13} = 9.2^\circ \quad (\text{A.8})$$

$$\theta_{23} = 38.5^\circ \quad (\text{A.9})$$

$$\Delta m^2 = 2.35 \times 10^{-3} \text{ eV}^2. \quad (\text{A.10})$$

These oscillation probabilities are first-order approximations to the full expressions. This is both for practical reasons – the second-order effects are poorly determined and are in fact what NOvA aims to measure – and because there is no requirement for the library have any particular distribution of events in it. The second order effects can pull the probabilities higher or lower, making this weighting a reasonable middle ground for the library. The library is also made devoid of intrinsic ν_e from the NuMI beam by setting that survival probability to zero. The overall prefactor on the $\nu_\mu \rightarrow \nu_e$ (signal) line relative to the background lines actually does not enter in practice since the signal, background, and π^0 -enriched background classes are scaled to have equal total weight in the library.

References

- [1] P. Adamson *et al.* (MINOS), Phys. Rev. Lett. **107**, 181802 (2011).
- [2] P. Adamson *et al.* (MINOS), Phys. Rev. Lett. **110**, 171801 (2013).
- [3] J. P. Ochoa, Ph.D. thesis, California Institute of Technology, FERMILAB-THESIS-2009-44 (2009).
- [4] R. Toner, Ph.D. thesis, University of Cambridge, FERMILAB-THESIS-2011-53 (2011).
- [5] D. S. Ayres *et al.* (NOvA), FERMILAB-DESIGN-2007-01 (2007); R. B. Patterson, for NOvA, Nucl. Phys. Proc. Suppl. **235-236**, 151 (2013).
- [6] K. Anderson *et al.*, FERMILAB-DESIGN-1998-01 (1998).
- [7] A. Ferrari, P. R. Sala, A. Fasso, and J. Ranft, CERN-2005-10 (2005); G. Battistoni *et al.*, AIP Conf. Proc. **896**, 31 (2007).
- [8] C. Andreopoulos, for GENIE, Acta Phys. Polon. B **40**, 2461 (2009).
- [9] S. Agostinelli *et al.*, Nucl. Instrum. Meth. A **506**, 250 (2003).
- [10] J. Uhlmann, Inform. Process. Lett. **40**, 4 (1991); P. N. Yianilos, Proceedings of the Fourth Annual ACM-SIAM Symposium on Discrete Algorithms, 311 (1993).