



Rauch<sup>bf</sup>, K. Rawlins<sup>c</sup>, I. C. Rea<sup>z</sup>, R. Reimann<sup>a</sup>, B. Relethford<sup>as</sup>, M. Renschler<sup>ad</sup>, G. Renzi<sup>m</sup>, E. Resconi<sup>z</sup>, W. Rhode<sup>v</sup>, M. Richman<sup>as</sup>, S. Robertson<sup>i</sup>, M. Rongen<sup>a</sup>, C. Rott<sup>ay</sup>, T. Ruhe<sup>v</sup>, D. Ryckbosch<sup>ab</sup>, D. Rysewyk<sup>w</sup>, I. Safa<sup>ak</sup>, S. E. Sanchez Herrera<sup>w</sup>, A. Sandrock<sup>v</sup>, J. Sandroos<sup>al</sup>, M. Santander<sup>ba</sup>, S. Sarkar<sup>ar</sup>, S. Sarkar<sup>x</sup>, K. Satalecka<sup>bf</sup>, M. Schaufel<sup>a</sup>, H. Schieler<sup>ad</sup>, P. Schlunder<sup>v</sup>, T. Schmidt<sup>r</sup>, A. Schneider<sup>ak</sup>, J. Schneider<sup>y</sup>, F. G. Schröder<sup>ad,ap</sup>, L. Schulte<sup>l</sup>, L. Schumacher<sup>a</sup>, S. Sclafani<sup>as</sup>, D. Seckel<sup>ap</sup>, S. Seunarine<sup>au</sup>, S. Shefali<sup>a</sup>, M. Silva<sup>ak</sup>, R. Snihur<sup>ak</sup>, J. Soedingrekso<sup>v</sup>, D. Soldin<sup>ap</sup>, S. Söldner-Rembold<sup>am</sup>, M. Song<sup>r</sup>, G. M. Spiczak<sup>au</sup>, C. Spiering<sup>bf</sup>, J. Stachurska<sup>bf</sup>, M. Stamatikos<sup>t</sup>, T. Stanev<sup>ap</sup>, R. Stein<sup>bf</sup>, P. Steinmüller<sup>ad</sup>, J. Stettner<sup>a</sup>, A. Steuer<sup>al</sup>, T. Stezelberger<sup>i</sup>, R. G. Stokstad<sup>i</sup>, A. Stöbl<sup>p</sup>, N. L. Strotjohann<sup>bf</sup>, T. Stürwald<sup>a</sup>, T. Stuttard<sup>u</sup>, G. W. Sullivan<sup>r</sup>, I. Taboada<sup>f</sup>, F. Tenholt<sup>k</sup>, S. Ter-Antonyan<sup>g</sup>, A. Terliuk<sup>bf</sup>, S. Tilav<sup>ap</sup>, K. Tollefson<sup>w</sup>, L. Tomankova<sup>k</sup>, C. Tönnis<sup>az</sup>, S. Toscano<sup>m</sup>, D. Tosi<sup>ak</sup>, A. Trettin<sup>bf</sup>, M. Tselengidou<sup>y</sup>, C. F. Tung<sup>f</sup>, A. Turcati<sup>z</sup>, R. Turcotte<sup>ad</sup>, C. F. Turley<sup>bc</sup>, B. Ty<sup>ak</sup>, E. Unger<sup>bd</sup>, M. A. Unland Elorrieta<sup>ao</sup>, M. Usner<sup>bf</sup>, J. Vandenbroucke<sup>ak</sup>, W. Van Driessche<sup>ab</sup>, D. van Eijk<sup>ak</sup>, N. van Eijndhoven<sup>n</sup>, J. van Santen<sup>bf</sup>, S. Verpoest<sup>ab</sup>, M. Vraeghe<sup>ab</sup>, C. Walck<sup>aw</sup>, A. Wallace<sup>b</sup>, M. Wallraff<sup>a</sup>, N. Wandkowsky<sup>ak</sup>, T. B. Watson<sup>d</sup>, C. Weaver<sup>x</sup>, A. Weindl<sup>ad</sup>, M. J. Weiss<sup>bc</sup>, J. Weldert<sup>al</sup>, C. Wendt<sup>ak</sup>, J. Werthebach<sup>ak</sup>, B. J. Whelan<sup>b</sup>, N. Whitehorn<sup>ag</sup>, K. Wiebe<sup>al</sup>, C. H. Wiebusch<sup>a</sup>, L. Wille<sup>ak</sup>, D. R. Williams<sup>ba</sup>, L. Wills<sup>as</sup>, M. Wolf<sup>z</sup>, J. Wood<sup>ak</sup>, T. R. Wood<sup>x</sup>, K. Woschnagg<sup>h</sup>, G. Wrede<sup>y</sup>, S. Wren<sup>am</sup>, D. L. Xu<sup>ak</sup>, X. W. Xu<sup>g</sup>, Y. Xu<sup>ax</sup>, J. P. Yanez<sup>x</sup>, G. Yodh<sup>ac</sup>, S. Yoshida<sup>p</sup>, T. Yuan<sup>ak</sup>, M. Zöcklein<sup>a</sup>

<sup>a</sup>*III. Physikalisches Institut, RWTH Aachen University, D-52056 Aachen, Germany*

<sup>b</sup>*Department of Physics, University of Adelaide, Adelaide, 5005, Australia*

<sup>c</sup>*Dept. of Physics and Astronomy, University of Alaska Anchorage, 3211 Providence Dr., Anchorage, AK 99508, USA*

<sup>d</sup>*Dept. of Physics, University of Texas at Arlington, 502 Yates St., Science Hall Rm 108, Box 19059, Arlington, TX 76019, USA*

<sup>e</sup>*CTSPS, Clark-Atlanta University, Atlanta, GA 30314, USA*

<sup>f</sup>*School of Physics and Center for Relativistic Astrophysics, Georgia Institute of Technology, Atlanta, GA 30332, USA*

<sup>g</sup>*Dept. of Physics, Southern University, Baton Rouge, LA 70813, USA*

<sup>h</sup>*Dept. of Physics, University of California, Berkeley, CA 94720, USA*

<sup>i</sup>*Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA*

<sup>j</sup>*Institut für Physik, Humboldt-Universität zu Berlin, D-12489 Berlin, Germany*

<sup>k</sup>*Fakultät für Physik & Astronomie, Ruhr-Universität Bochum, D-44780 Bochum, Germany*

<sup>l</sup>*Physikalisches Institut, Universität Bonn, Nussallee 12, D-53115 Bonn, Germany*

<sup>m</sup>*Université Libre de Bruxelles, Science Faculty CP230, B-1050 Brussels, Belgium*

<sup>n</sup>*Vrije Universiteit Brussel (VUB), Dienst ELEM, B-1050 Brussels, Belgium*

<sup>o</sup>*Dept. of Physics, Massachusetts Institute of Technology, Cambridge, MA 02139, USA*

<sup>p</sup>*Dept. of Physics and Institute for Global Prominent Research, Chiba University, Chiba 263-8522, Japan*

<sup>q</sup>*Dept. of Physics and Astronomy, University of Canterbury, Private Bag 4800, Christchurch, New Zealand*

<sup>r</sup>*Dept. of Physics, University of Maryland, College Park, MD 20742, USA*

<sup>s</sup>*Dept. of Astronomy, Ohio State University, Columbus, OH 43210, USA*

<sup>t</sup>*Dept. of Physics and Center for Cosmology and Astro-Particle Physics, Ohio State University, Columbus, OH 43210, USA*

<sup>u</sup>*Niels Bohr Institute, University of Copenhagen, DK-2100 Copenhagen, Denmark*

<sup>v</sup>*Dept. of Physics, TU Dortmund University, D-44221 Dortmund, Germany*

<sup>w</sup>*Dept. of Physics and Astronomy, Michigan State University, East Lansing, MI 48824, USA*

<sup>x</sup>*Dept. of Physics, University of Alberta, Edmonton, Alberta, Canada T6G 2E1*

- <sup>y</sup>*Erlangen Centre for Astroparticle Physics, Friedrich-Alexander-Universität Erlangen-Nürnberg, D-91058 Erlangen, Germany*
- <sup>z</sup>*Physik-department, Technische Universität München, D-85748 Garching, Germany*
- <sup>aa</sup>*Département de physique nucléaire et corpusculaire, Université de Genève, CH-1211 Genève, Switzerland*
- <sup>ab</sup>*Dept. of Physics and Astronomy, University of Gent, B-9000 Gent, Belgium*
- <sup>ac</sup>*Dept. of Physics and Astronomy, University of California, Irvine, CA 92697, USA*
- <sup>ad</sup>*Karlsruhe Institute of Technology, Institut für Kernphysik, D-76021 Karlsruhe, Germany*
- <sup>ae</sup>*Dept. of Physics and Astronomy, University of Kansas, Lawrence, KS 66045, USA*
- <sup>af</sup>*SNOLAB, 1039 Regional Road 24, Creighton Mine 9, Lively, ON, Canada P3Y 1N2*
- <sup>ag</sup>*Department of Physics and Astronomy, UCLA, Los Angeles, CA 90095, USA*
- <sup>ah</sup>*School of Physics and Astronomy, Queen Mary University of London, London E1 4NS, United Kingdom*
- <sup>ai</sup>*Department of Physics, Mercer University, Macon, GA 31207-0001, USA*
- <sup>aj</sup>*Dept. of Astronomy, University of Wisconsin, Madison, WI 53706, USA*
- <sup>ak</sup>*Dept. of Physics and Wisconsin IceCube Particle Astrophysics Center, University of Wisconsin, Madison, WI 53706, USA*
- <sup>al</sup>*Institute of Physics, University of Mainz, Staudinger Weg 7, D-55099 Mainz, Germany*
- <sup>am</sup>*School of Physics and Astronomy, The University of Manchester, Oxford Road, Manchester, M13 9PL, United Kingdom*
- <sup>an</sup>*Department of Physics, Marquette University, Milwaukee, WI, 53201, USA*
- <sup>ao</sup>*Institut für Kernphysik, Westfälische Wilhelms-Universität Münster, D-48149 Münster, Germany*
- <sup>ap</sup>*Bartol Research Institute and Dept. of Physics and Astronomy, University of Delaware, Newark, DE 19716, USA*
- <sup>aq</sup>*Dept. of Physics, Yale University, New Haven, CT 06520, USA*
- <sup>ar</sup>*Dept. of Physics, University of Oxford, Parks Road, Oxford OX1 3PU, UK*
- <sup>as</sup>*Dept. of Physics, Drexel University, 3141 Chestnut Street, Philadelphia, PA 19104, USA*
- <sup>at</sup>*Physics Department, South Dakota School of Mines and Technology, Rapid City, SD 57701, USA*
- <sup>au</sup>*Dept. of Physics, University of Wisconsin, River Falls, WI 54022, USA*
- <sup>av</sup>*Dept. of Physics and Astronomy, University of Rochester, Rochester, NY 14627, USA*
- <sup>aw</sup>*Oskar Klein Centre and Dept. of Physics, Stockholm University, SE-10691 Stockholm, Sweden*
- <sup>ax</sup>*Dept. of Physics and Astronomy, Stony Brook University, Stony Brook, NY 11794-3800, USA*
- <sup>ay</sup>*Dept. of Physics, Sungkyunkwan University, Suwon 16419, Korea*
- <sup>az</sup>*Institute of Basic Science, Sungkyunkwan University, Suwon 16419, Korea*
- <sup>ba</sup>*Dept. of Physics and Astronomy, University of Alabama, Tuscaloosa, AL 35487, USA*
- <sup>bb</sup>*Dept. of Astronomy and Astrophysics, Pennsylvania State University, University Park, PA 16802, USA*
- <sup>bc</sup>*Dept. of Physics, Pennsylvania State University, University Park, PA 16802, USA*
- <sup>bd</sup>*Dept. of Physics and Astronomy, Uppsala University, Box 516, S-75120 Uppsala, Sweden*
- <sup>be</sup>*Dept. of Physics, University of Wuppertal, D-42119 Wuppertal, Germany*
- <sup>bf</sup>*DESY, D-15738 Zeuthen, Germany*

---

## Abstract

The current and upcoming generation of Very Large Volume Neutrino Telescopes—collecting unprecedented quantities of neutrino events—can be used to explore subtle effects in oscillation physics, such as (but not restricted to) the neutrino mass ordering. The sensitivity of an experiment to these effects can be estimated from Monte Carlo simulations. With the high number of events that will be collected, there is a trade-off between the computational expense of running such simulations and the inherent statistical uncertainty in the determined values. In such a scenario, it becomes impractical to produce and use

adequately-sized sets of simulated events with traditional methods, such as Monte Carlo weighting. In this work we present a staged approach to the generation of binned event distributions in order to overcome these challenges. By combining multiple integration and smoothing techniques which address limited statistics from simulation it arrives at reliable analysis results using modest computational resources.

*Keywords:* Data Analysis, Monte Carlo, MC, Statistics, Smoothing, KDE, Neutrino, Neutrino Mass Ordering, Detector, VLV $\nu$ T

---

## 1. Introduction

By virtue of their multi-megaton effective mass paired with the magnitude of the atmospheric neutrino flux, the next generation of Very Large Volume Neutrino Telescopes (VLV $\nu$ Ts) dedicated to neutrino oscillation physics, such as the IceCube Upgrade, PINGU, and ORCA [1, 2, 3, 4], will record tens of thousands of GeV-scale neutrino interactions. These large-scale water or ice Cherenkov detectors do not have the ability to unambiguously distinguish between neutrino flavors and interaction types on an event-by-event basis. Even so, their high statistics data samples can be used to explore effects that are small compared to the background, such as the tau neutrino appearance rate, the ordering of the neutrino mass eigenstates (NMO), or potential neutrino physics beyond the Standard Model.

All such physics analyses are carried out by comparing the observed event distributions with predictions (hereafter referred to as *templates*) obtained from Monte Carlo (MC) simulations. The physical phenomena listed above will appear as statistical (in)compatibilities of templates with differences in event counts as small as a few percent. An inherent problem when trying to quantify these deviations in high-statistics data sets is that the templates must be described with an accuracy better than the magnitude of the effect being investigated. A limiting factor to the accuracy is the amount of MC simulation available, which is in turn constrained by the availability of computing resources. This particularly applies during the design optimization phase of a planned experiment, which entails performance assessments of multiple detector variants.

With an adequate machinery at hand to produce templates, extracting the relevant physical and systematic parameters typically proceeds via maximizing the likelihood of obtaining the observed data under a given hypothesis. A common feature to all statistical methods is that the templates need to be generated for a multitude of parameter combinations, often thousands or even millions. This process needs to be accurate, but also fast, which typically prohibits the reproduction of the full MC sample for each template.

---

\*analysis@icecube.wisc.edu

<sup>1</sup>also at Università di Padova, I-35131 Padova, Italy

<sup>2</sup>also at National Research Nuclear University, Moscow Engineering Physics Institute (MEPhI), Moscow 115409, Russia

<sup>3</sup>Earthquake Research Institute, University of Tokyo, Bunkyo, Tokyo 113-0032, Japan

In this article, we present an approach that allows for the fast creation of accurate templates even from MC sets that are several orders of magnitude smaller than those necessary when using simpler methods. An alternative approach that does not remove template inaccuracies but rather mitigates their impact on statistical inference is the inclusion of the inherent MC uncertainty in the fit statistic; recent overviews can be found in [5, 6].

Our approach was used to calculate the expected sensitivities for atmospheric neutrino oscillation analyses with the proposed PINGU experiment [2, 3], and a similar approach was taken in low-energy sensitivity studies for the KM3NeT design [4]. Throughout this article, we will use the NMO analysis for a generic VLV $\nu$ T as an example to illustrate our methods, though it is applicable in a wider range of atmospheric neutrino oscillation analyses, and, in parts and with limitations, to other experiments. Section 2 details the computational challenge at hand, followed by a brief introduction of the example NMO analysis in Section 3. Our approach to overcome this challenge is presented in Section 4 and Section 5, followed by a discussion of the validity of the approach in Section 6. The performance is compared to other typical analysis methods in Section 7, while the computational burden is discussed in Section 8. Section 9 concludes with a brief summary of the article. Finally, in Appendix A we provide details about the VLV $\nu$ T toy model that we use to benchmark the performance of all considered analysis approaches.

## 2. Computational Challenge

The statistical comparison between experimental data and parametric or MC-based predictions allows inference of the values of physics parameters under study. It typically proceeds via a likelihood analysis. We first discuss its most general concepts and variants, then detail the computational requirements on MC generation, and finally outline two standard methods of mitigating these computational burdens.

### 2.1. Likelihood Analysis

Different types of likelihood analyses in particle physics share common features<sup>4</sup>. An experiment records data which are used to reconstruct any observables expected to carry the imprint of the physical phenomenon under study. A selection (triggering, filtering, etc.) is applied in order to enhance the sought signal. Before performing statistical inference, we need a theoretical model of the observable distributions to compare to the data. Often this includes complicated processes like particle interactions and detector response that require the use of MC methods. Hence, not only the data, but also the model is subject to statistical fluctuations. However, once an appropriate amount of MC events is available, the data  $x_i$  can be compared to templates—theoretical distributions—for different physics parameter values  $\boldsymbol{\theta}$  via a likelihood function,  $L(x_1, x_2, \dots, x_n|\boldsymbol{\theta}) = \prod_i P(x_i|\boldsymbol{\theta})$ , where  $P(x_i|\boldsymbol{\theta})$  is the probability to observe the data  $x_i$  assuming that  $\boldsymbol{\theta}$  corresponds to given values of the physics parameters<sup>5</sup>.

---

<sup>4</sup>See, for example, [7] for a more complete overview.

<sup>5</sup>If the total number of measurements,  $n$ , is also a random quantity, the likelihood function can be extended to include the distribution of  $n$  [8].

The goal is (in the frequentist picture) to find the maximum likelihood estimators (MLEs)  $\hat{\boldsymbol{\theta}}$ , i.e., the parameter values which maximize  $L$ .

The methods presented in this paper depend on a likelihood function applied to binned data. One usually employs either a Poisson likelihood or a  $\chi^2$  approximation thereof (see for example [8]); our methods are independent of this choice, but we use the latter in the example presented in this article. Binning the data hides physics signatures smaller than the bin size and thus introduces a loss in sensitivity. This can be mitigated by reducing bin sizes, but smaller bins come at the cost of reduced—and possibly insufficient—MC statistics in each bin.

Apart from the physics parameters of interest, a model often comes with nuisance parameters that are also included in the likelihood function. This further increases the dimensionality of the MLE search, which relies on numerical routines for multidimensional optimization problems. For the NMO studies, we use the L-BFGS-B algorithm [9] in a  $D = 8$  dimensional parameter space (see Table 1). The number of steps necessary for the optimization to converge depends on the particular analysis and model being used (i.e., the details of the resulting likelihood landscape); in the case of our toy example, an average of  $\sim 10^3$  templates (one per realization of  $\boldsymbol{\theta}$ ) were needed to converge.

## 2.2. Template and MC Generation Requirements

The problems associated with generating such a large number of templates are exacerbated when estimating the median sensitivity of an experiment. The above process needs to be applied to an ensemble of random toy MC pseudo-experiments<sup>6</sup> of size  $N_p$ . The comparison of test statistic distributions  $\mathcal{T}$  (see Section 3 for details) can be used to estimate a significance value  $n_\sigma$  at which one hypothesis is preferred over the alternative. If  $\mathcal{T}$  is Gaussian distributed<sup>7</sup>, the uncertainty  $\Delta n_\sigma$  to which  $n_\sigma$  can be determined depends upon the number of pseudo-experiments  $N_p$  as (see Appendix B for details):

$$\Delta n_\sigma = \frac{1}{\sqrt{N_p}} \sqrt{\frac{n_\sigma^2}{2} + 2}. \quad (1)$$

With an absolute uncertainty  $\Delta n_\sigma$  at the 1% level, determining the sensitivity of an experiment at a confidence level of 99.7% (corresponding to  $n_\sigma = 3$ ) requires  $\mathcal{O}(10^4)$  pseudo-experiments.

Finally, the event count expectations,  $\boldsymbol{\mu}$ , for all bins in the templates must be determined at the same level as the physics effects being investigated, which requires at least  $\frac{1}{(1\%)^2} = 10^4$  MC events per bin to study sub-percent variations arising in a comparison of the two NMO realizations. At the same time, the number of bins used in any histograms must be

---

<sup>6</sup>Each pseudo-experiment corresponds to a statistical fluctuation of the expected experimental outcome as predicted by MC events. For certain problems, the generation of pseudo-experiments can be skipped by applying the *Asimov* approximation [10, 3].

<sup>7</sup>While not a prediction from the model, a near-Gaussian distribution of the test statistic is observed in most NMO studies [3, 4, 11].

commensurate with the experimental resolution and the feature size of the effect under study. In the example case, at least  $\mathcal{O}(10^3)$  bins are required to resolve the distinct features of the NMO signature; otherwise the analysis cannot exploit the full potential of the experiment.

Therefore, the brute-force approach to our example case requires a very large number of neutrino events to be simulated:  $\mathcal{O}(10^7)$  events for each of  $\mathcal{O}(10^3)$  values of  $\theta$  probed during the optimization process for each of  $\mathcal{O}(10^4)$  pseudo-experiments—a grand total of  $\mathcal{O}(10^{14})$  events. Even if the time to simulate and reconstruct a single event is 1 s (a very optimistic estimate for our experiment), full fits to all pseudo-experiments under the two ordering hypotheses would require  $\mathcal{O}(10^{10})$  CPU-core-hours—i.e., a single analysis would keep  $10^5$  CPU cores busy for 30 years<sup>8</sup>—a restriction clearly prohibitive to performing any study. Various state-of-the-art methods are employed to mitigate the high computational costs. In the remainder of this section, we briefly present the main ideas behind these methods and give a conceptual introduction to how they are embedded in the approach we introduce in this article.

### 2.3. Weighting and Smoothing

The standard event-by-event MC weighting technique avoids repeated simulation and reconstruction of events every time a value of a nuisance parameter is changed. This is possible, first, because the physics processes of neutrino production in the atmosphere (flux), their propagation involving flavor oscillation, and their detection and reconstruction are independent. Each of these processes, therefore, can be treated separately.

For a process that has an a priori known parametric form (the parameter values of which are not necessarily known), the outcome of that process can be predicted by directly evaluating the parametrization at a set of input values. In our case, both the neutrino flux prediction and flavor oscillations fall into this category. The second category of processes are those that require MC simulation. Predictions of the detection and reconstruction of neutrinos fall into this category because we do not have a complete characterization of the detector’s response.

This leads to the standard event-by-event reweighting scheme, which estimates the expected final-level event counts due to all processes by simulating a set of MC neutrinos (capturing the effects of detection and reconstruction), assigning to each a weight derived from flux and oscillation calculations, and binning the events’ weights in some set of observable dimensions, as illustrated in the top row of Figure 1.

In detail: Each MC neutrino—generated with a flavor  $\beta$  and a set of true observables  $\theta_\nu^{\text{true}}$ —is assigned a posteriori the weight  $w_\beta$  corresponding to the sum over the atmospheric fluxes  $\Phi_\alpha(\theta_{\text{flux}}; \theta_\nu^{\text{true}})$  of all initial flavors  $\alpha$  including the probabilities  $P_{\alpha \rightarrow \beta}^{\text{osc}}(\theta_{\text{osc}}; \theta_\nu^{\text{true}})$  to oscillate into a neutrino of the flavor  $\beta$ :

$$w_\beta \propto \sum_{\alpha} \Phi_{\alpha}(\theta_{\text{flux}}; \theta_{\nu}^{\text{true}}) \times P_{\alpha \rightarrow \beta}^{\text{osc}}(\theta_{\text{osc}}; \theta_{\nu}^{\text{true}}).$$

---

<sup>8</sup>Here we make the assumption that the algorithm can be parallelized perfectly.

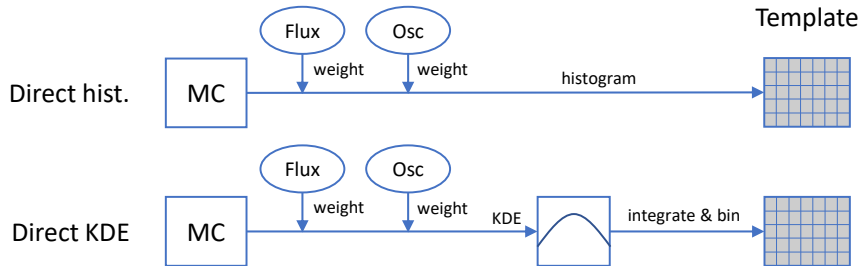


Figure 1: Operating principles of direct histogramming (top row) and direct KDE (bottom row), which both follow the same weighting scheme for MC events but arrive at the template differently, as explained in the text.

In the above,  $\theta_{\text{flux}}$  and  $\theta_{\text{osc}}$  are nuisance parameters affecting neutrino fluxes and oscillation probabilities. For a given realization of non-parametric nuisance parameters  $\theta_{\text{det}}$  affecting the detector response, applying the detector response simulation (including event reconstruction, classification, etc.) to each incident neutrino results in a set of reconstructed observables  $\theta_{\nu}^{\text{reco}}$ , whose distribution can be compared to real data. In practice, techniques dealing with the discrete nature of detector nuisance parameters may be required. Here, however, we consider only a single realization of the detector parameters ( $\theta_{\text{det}}$  fixed)—a simplification without any loss of the general applicability of the methods discussed.

Since the process of oscillation is decoupled from the detector simulation, only a single MC set is required to generate the templates for the different hypotheses under test (e.g., the two mass orderings); only the weights  $w_{\beta}$  must be recomputed. This eliminates statistical fluctuations between the otherwise disjoint MC samples. However, even with a single MC set, an undersampling of the phase space of the model can result in a bias.

Binning the weights in (a relevant subset of)  $\theta_{\nu}^{\text{reco}}$  corresponds to performing MC integration of the experiment’s event distribution. While the convergence rate of this approach does not depend on the dimensionality of the integral, errors of the estimates scale as  $1/\sqrt{N}$ , where  $N$  is the number of MC events that fall in a bin.

As it is often infeasible to generate enough MC events to obtain sufficient accuracy in the MC integration process, smoothing of the final event distributions is a common practice. This, however, can be computationally slow and can introduce artificial features which may incorrectly reduce or enhance the signal. One such smoothing technique is kernel density estimation (KDE) [12]. Specifically, we apply adaptive bandwidth KDE directly to the weighted MC to compare a state-of-the-art version of this method to the methods we introduce in this paper in Section 4. Here, a Gaussian kernel with a width calculated as described in [13] is centered at each MC event’s reconstruction information. A weighted sum over the kernels of all events then delivers the smoothed distribution as shown in the bottom row of Figure 1, which will be compared to the distribution our method yields.

Shortcomings of the direct application of the two techniques discussed above—the first is the weighting method alone (labeled *direct histogramming*), while the second applies additional smoothing using adaptive kernel density estimates (labeled *direct KDE*)—can be



overcome using the *staged approach*. Before providing an overview of the staged approach in Section 4, we briefly introduce the key points of the example NMO analysis used to illustrate the benefits of the approach with respect to the standard techniques.

### 3. NMO Analysis

The observation of neutrino oscillations and the demonstration of the neutrinos' non-zero masses [14, 15] represented a major step forward in the field of particle physics. While current experimental techniques have not yet allowed for a direct measurement of the tiny masses, the magnitudes of their relative differences (mass splittings) are well known.

The ordering of these neutrino mass states (neutrino mass ordering, NMO) presents a difficult challenge. A powerful method to determine this ordering is the observation of matter effects on neutrinos. Owing to the high electron density of the Sun, observations of solar neutrinos have shown the second mass state to be heavier than the first [16]. It remains an open question, however, whether the third state is the most or least massive. The former scenario is referred to as the normal ordering (NO), while the second is called inverted ordering (IO). There is currently no experimental evidence decisively excluding either of the two scenarios [17, 18, 19, 20].

The study of oscillations of atmospheric neutrinos provides a promising route toward a decisive measurement of the NMO [21, 2, 3, 4]. The path length (or *baseline*) varies between 20 km for vertically downward going and 12 700 km for straight upward going atmospheric neutrinos, with the latter crossing the full diameter of the Earth. With energies ranging from MeV up to the TeV scale, combinations of baselines and energies varying over several orders of magnitude are probed. For the longest baseline, the very pronounced first oscillation maximum of muon neutrinos occurs at a neutrino energy of around 25 GeV, followed by subsequent maxima at lower energies.

The electron neutrinos' coupling to electrons (coherent forward scattering) in the Earth creates an effective matter potential which leads to resonant behavior of the transition probabilities at energies around 5 GeV, known as matter resonances [22, 23, 24]. Furthermore, the Earth's specific density profile encountered by the neutrinos can also parametrically enhance their oscillations [25]. This enhancement with respect to oscillations proceeding in vacuum occurs for neutrinos if the NMO is normal, otherwise for anti-neutrinos.

The NMO measurement potential of VLV $\nu$ Ts is based on this asymmetry. Two major aspects are obstructive, however. The first is the inability of VLV $\nu$ Ts to differentiate between neutrinos and anti-neutrinos. This reduces the effect to the respective difference in atmospheric fluxes and interaction cross sections. Energy and directional resolutions of the experiment present the second hurdle. Both are typically prohibitive to resolving the fast variations of the oscillation pattern at the relevant energies. As a consequence, the observable effect is reduced to at most a few percent over the relevant energy and zenith range (see Figure 2), requiring neutrino telescopes with effective masses on the order of megatons to achieve sufficient event statistics.

Proponents of various VLV $\nu$ Ts in ice and water have performed studies confirming this idea, finding that a  $> 3\sigma$  (median) sensitivity to the NMO can be achieved within five

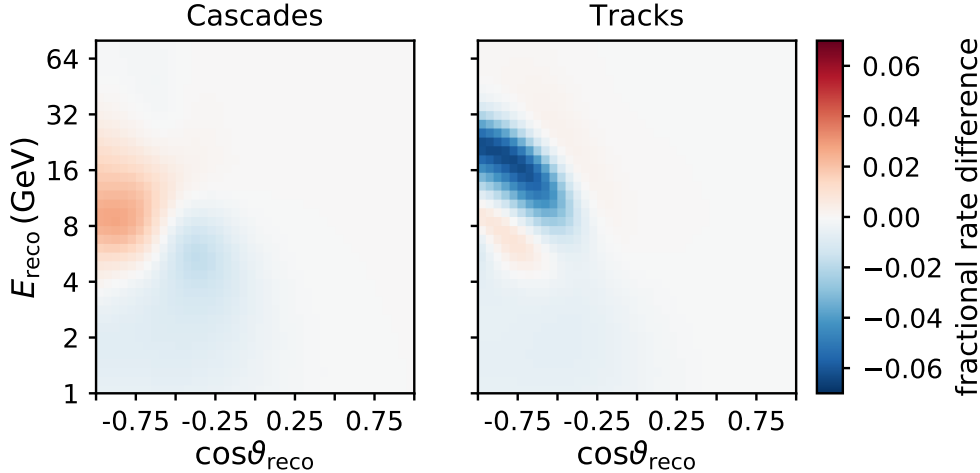


Figure 2: Expected fractional event rate difference between nominal NO and IO inputs (from Table 1) for the toy model. Cascades are shown on the left, tracks on the right.

years of exposure time even in less favorable regions of the neutrino oscillation parameter space [2, 4, 26].

As the oscillation probabilities directly depend on neutrino energy  $E_{\text{true}}$ , oscillation baseline ( $\propto \cos \vartheta_{\text{true}}$ ), and flavor, we split our data into bins of  $\log_{10} E_{\text{reco}}$ ,  $\cos \vartheta_{\text{reco}}$ , and event class<sup>9</sup>. It is important to choose a binning fine enough to resolve the NMO signature, while coarse enough to retain a sufficient amount of MC statistics per bin, as motivated in Section 2. We have found the division into  $(40 \times 40 \times 2)$  bins to be suitable, covering a range of  $E_{\text{reco}}$  from 1 GeV to 80 GeV, the whole sky ( $\cos \vartheta_{\text{reco}}$  from  $-1$  to  $1$ ), and the two event classes of *cascades* and *tracks*. Using this binning, for our toy detector introduced in Appendix A, Figure 2 shows the expected fractional event rate difference  $(R_{\text{NO}} - R_{\text{IO}})/R_{\text{NO}}$ , where  $R_{\text{NO(IO)}}$  is the expected event rate for true NO (IO), based on the two sets of nominal model parameter values given in Table 1.

As the most powerful test statistic for distinguishing two simple hypotheses [27], the *logarithm of the likelihood ratio*

$$\mathcal{T} = -2 \ln \left( \frac{\max_{\boldsymbol{\theta} \in \text{NO}} L(\mathbf{n} | \boldsymbol{\mu}(\boldsymbol{\theta}))}{\max_{\boldsymbol{\theta} \in \text{IO}} L(\mathbf{n} | \boldsymbol{\mu}(\boldsymbol{\theta}))} \right). \quad (2)$$

is also useful in assessing the ability of an experiment to discriminate between the two (composite) NMO hypotheses at a given confidence level. It is representative of the degree at which observing the data  $\mathbf{n}$  under the NO hypothesis is favored over observing it under the alternate IO hypothesis. The observed spectrum at the detector,  $\mathbf{n}$ , however, is a convolution

<sup>9</sup>The use of the subscript “true” is used to specify the true variables of the neutrinos and to distinguish these from the reconstructed variables, denoted with a subscript “reco”, which will be introduced in Section 4.1.

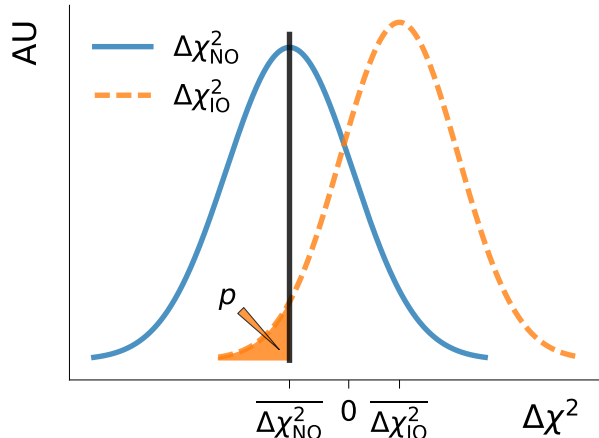


Figure 3: Example distributions of Equation (3). The distribution on the left (solid line) represents the case of NO pseudo-data, while the distribution on the right (dashed) is obtained when the pseudo-data is taken from the IO. Here,  $1 - p$  corresponds to the confidence level at which the IO is correctly rejected with a probability of 50%.

of the atmospheric neutrino flux, the effects of neutrino oscillations that bear the NMO signature, the neutrino interaction and detection processes, and the event reconstruction and classification procedure. Each one of these effects is accompanied by systematic uncertainties. As their impact on the predicted spectrum  $\mu$  is modeled, the systematic uncertainties directly feed in to the likelihood  $L$  of the observation.

For this study, we limit ourselves to a simplified treatment using  $\chi^2$  statistics and the *Asimov* dataset. In this approach, the projected median sensitivity is calculated from the average experimental outcomes under the two possible NMO hypotheses, as opposed to performing extensive ensemble tests with randomly fluctuated pseudo-experiments. The log-likelihood expression is a simple  $\chi^2$ , and Equation (2) can be rewritten as the difference

$$\Delta\chi^2 = \chi_{\text{NO}}^2 - \chi_{\text{IO}}^2. \quad (3)$$

Here,  $\chi_{\text{NO}}^2$  is the minimum  $\chi^2$  between model predictions and data, with all nuisance parameters profiled out using NO priors ( $\chi_{\text{IO}}^2$  follows analogously).

An illustration of example distributions of (3) for the two different NMO hypotheses is shown in Figure 3. The goal is to obtain a p-value  $p$  which quantifies the statistical compatibility between the hypothesis that is tested and the one assumed to be true. In the ensemble approach, the two distributions would need to be built up by fitting pseudo-experiments. In the Asimov approach, however, certain assumptions about the distribution of (3) allow adopting the expression  $\sqrt{|\Delta\chi^2|}$  as a sensitivity proxy [11], determining the significance at which the wrong ordering can be excluded without the need for pseudo-experiments.

For the profiling of the nuisance parameters (any free model parameters), a numerical algorithm minimizes the  $\chi^2$  metric. Whenever external constraints are applied to such parameters, we add those to the  $\chi^2$  value as penalty terms (priors). While the presence of

Parameter	Nominal value		Prior
	NO	IO	
$\nu_e/\nu_\mu$ flux ratio	1.0	1.0	$\pm 0.03$
$\nu/\bar{\nu}$ flux ratio	1.0	1.0	$\pm 0.1$
Spectral index shift	0.0	0.0	$\pm 0.1$
Energy scale	1.0	1.0	$\pm 0.1$
Overall normalization	1.0	1.0	$\pm 0.1$
$\theta_{13}$ ( $^\circ$ )	8.5	8.5	$\pm 0.2$ [28, 29]
$\theta_{23}$ ( $^\circ$ )	42.3	49.5	non-Gaussian [28, 29]
$\Delta m_{31}^2$ ( $\text{eV}^2$ )	0.00246	-0.00237	$\pm 4.75 \times 10^{-5}$ [28, 29]

Table 1: Summary of model parameters in the example NMO analysis, including their nominal values for the two NMO hypotheses and Gaussian  $\pm 1\sigma$  bounds used as external constraints (priors). The first three parameters are applied to atmospheric neutrino flux predictions from [30], following the procedure laid out in Section 5.1. The values for the three oscillation parameters are based on a recent global fit [28, 29].

these penalty terms is meant to illustrate a typical approach to problems of this sort, their sizes do not follow any precise physical motivation. Table 1 gives an overview of all used model parameters, their nominal values for NO and IO, and priors (where applied).

## 4. Overview of the Staged Approach

The method to obtain templates we describe in this article is divided into four independent parts, referred to as *stages*. The four stages (flux, oscillation, detection, and reconstruction) and how they are used to obtain event templates are summarized in this section, while more technical descriptions of each stage follow in Section 5.

### 4.1. Stages

Templates for our example case of an NMO analysis using a VLV $\nu$ T are produced efficiently and accurately using the following four stages.

Each stage represents a collection of related physical effects. Beginning with the flux computed by the initial stage, each subsequent stage applies a transformation to the output of the previous stage.

1. **Flux** The expected unoscillated atmospheric neutrino fluxes are taken from an external model [30]. Flux values from this model are provided in the form of tables with discrete steps in both neutrino energy,  $E_{\text{true}}$ , and direction, here the cosine of the zenith angle,  $\cos \vartheta_{\text{true}}$ . Therefore, an interpolation must be performed for values between those tabulated. Crucially, these tables give the integrated flux across the bins, which does not necessarily coincide with the flux value at the bin center. Accordingly, we use an integral-preserving (IP) interpolation. In general, atmospheric flux models require external inputs including primary cosmic ray measurements, atmospheric density models, and hadronic interaction measurements. Many associated uncertainties are known [31, 32] and need to be included as nuisance parameters in an analysis.

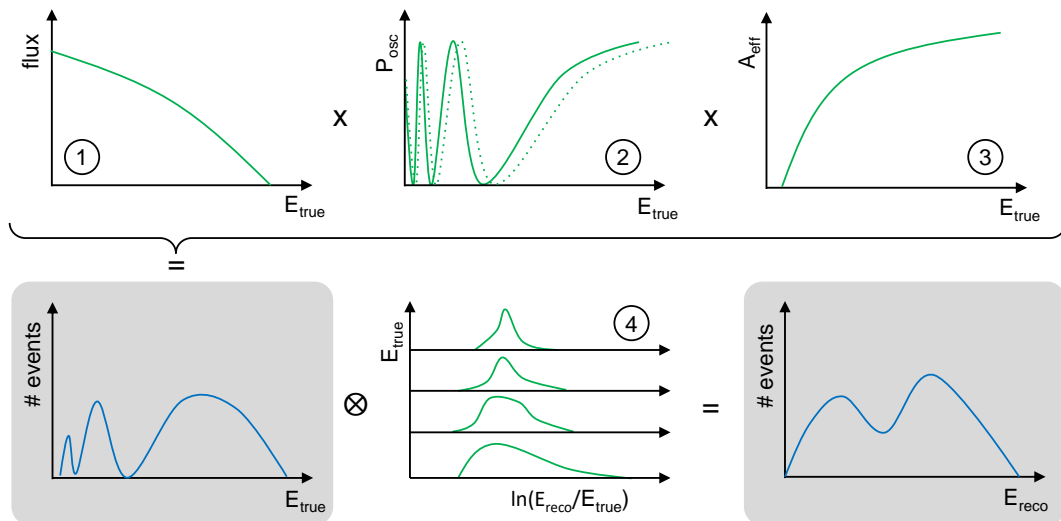


Figure 4: Illustration of the staged approach for obtaining event templates, here for simplicity using a characterization in one dimension (energy) only. Steps 1, 2, and 3 are in true energy ( $E_{\text{true}}$ ); the product of these yields the expected event distribution (lower left). Smearing this spectrum with energy-dependent energy resolution functions (step 4) gives the reconstructed event rate spectrum (lower right). Note that the dotted green line in step 2 shows a hypothetical change of oscillation parameters, affecting only stage 2. Smoothing can now directly be applied to the distributions in steps 3 and 4, instead of the fully weighted MC as in the direct KDE method.

2. **Oscillation** Flavor oscillations of neutrinos traversing the Earth modify the flavor content of the original flux in a manner that depends on the energies and path lengths (derived from the direction) of the neutrinos. Additional intrinsic neutrino properties determine the standard flavor oscillation probabilities: three mixing angles and two independent mass-squared splittings, as well as a possible non-zero CP-violating phase. In addition, matter effects induce modifications in the flavor transition probabilities compared to vacuum [23, 22, 33], which makes up the basis of the NMO measurement capability of VLV $\nu$ Ts. In [33], the authors present an analytical expression for the neutrino flavor transition amplitude in a layer of uniform-density matter, which in turn was later implemented in, for example, the `Prob3++` software [34]. Here, the Earth density profile [35] is approximated by a finite number of homogeneous layers and the total transition amplitude is represented by a matrix product of the amplitudes in the individual layers. The main challenge for this stage, which in contrast to the other stages does not require any MC simulation, is to keep the burden of these computationally expensive calculations to a minimum, while retaining sufficient accuracy in the modeling of the neutrinos' propagation.
3. **Detection** The number of observed events is determined by the (oscillated) flux as well as a quantity known as the *effective area* (alternatively, the effective mass)<sup>10</sup>. This

<sup>10</sup>In contrast, high energy physics experiments often calculate an acceptance instead, which is also based on simulation.

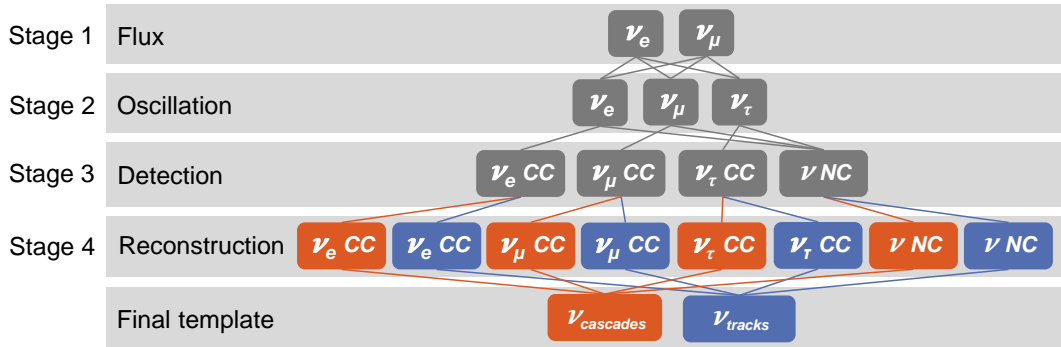


Figure 5: Flow of neutrino flavors and interaction types through the stages, here shown for neutrinos only (with an analogue counterpart for anti-neutrinos). Neutral current events of all flavors are indistinguishable and can therefore be conveniently added together. The reconstruction stage not only maps from  $(E_{\text{true}} \times \cos \vartheta_{\text{true}})$ -space to  $(E_{\text{reco}} \times \cos \vartheta_{\text{reco}})$ -space, but also classifies the events into the cascade and track categories, indicated by the orange and blue color, respectively.

incorporates the probability that a given neutrino interacts within the detector, is detected, and passes the given data selection criteria. We obtain the eight effective areas ( $\nu_{e,\mu,\tau}$  &  $\bar{\nu}_{e,\mu,\tau}$  charged current (CC) and  $\nu$  &  $\bar{\nu}$  neutral current (NC) interactions) from simulated MC events that are run through the same selection criteria as the real data. In general, each of these effective areas will depend on the energy and arrival direction of the neutrinos. Depending on the detector geometry, certain symmetries can be exploited to reduce the number of parameters on which the effective areas depend. Here we assume azimuthal symmetry and therefore only extract effective areas as a function of  $E_{\text{true}}$  and  $\cos \vartheta_{\text{true}}$ .

4. **Reconstruction** The process referred to as *reconstruction* translates the raw signals recorded by a detector into estimates of the physical properties of events. Uncertainties in these estimates manifest as statistical fluctuations, with respect to the true properties, which can be described by probability density functions we refer to as *resolution functions*. We estimate the resolution functions from the same MC events as used in the detection stage, for which we know the true energy, zenith angle, and interaction type on an event-by-event level. The reconstruction stage uses these estimated resolution functions to build smearing kernels (ensembles of resolution functions) that map the event rates from the space of true variables into the space of reconstructed observables. Additionally—since most VLV $\nu$ Ts cannot exactly distinguish the different neutrino flavors and interaction types—the events are classified by their signature in the detector. Here, event classes are *tracks* and *cascades*, based on whether the event seems to contain the expected signature of a starting muon track. This process will separate  $\nu_\mu$  CC and  $\bar{\nu}_\mu$  CC interactions from all others, albeit with limited efficiency and purity. For the example NMO analysis, three observables are needed: the primary neutrino’s reconstructed energy ( $E_{\text{reco}}$ ), zenith angle ( $\vartheta_{\text{reco}}$ ), and event classification.

Note that there is no universal prescription for identifying the set of stages appropriate

for any given physics analysis or detector. Instead, stages are chosen to exploit valid simplifications for the task at hand. For example, atmospheric neutrino flux and oscillation calculations depend on readily available tabulated spectra and analytic formulae, respectively. Cosmic ray observatories or high energy particle colliders, by contrast, might require complex stages to describe particle showers, which in turn might depend on high-dimensional, analysis-specific tables. Any physics scenario resulting in multi-particle final states adds further complexity.

In essence, the specific problem and analysis at hand determine to which extent MC sampling is necessary and whether the staged approach is applicable. If the latter is indeed the case, care must be taken concerning the choice of appropriate stages and their specific implementations. In the remainder of this article, we study in detail the staged approach we have found particularly effective for an NMO analysis using a VLV $\nu$ T.

#### 4.2. Template Generation

In order to produce the final-level event templates that are ultimately compared to the data, the four stages are combined as depicted in Figure 4: integration of the product of the first three stages (flux, oscillation probability, and effective area) over  $E_{\text{true}}$  and  $\cos\vartheta_{\text{true}}$  yields the event rate in the space of true variables. The event rate in the space of reconstructed observables is then obtained by a convolution of the true event rate with the reconstruction resolution functions. Finally, multiplication by detector exposure time results in an event *count*, which can be compared directly to observed data or different templates<sup>11</sup>.

Throughout the stages, different combinations of neutrino flavor and interaction type (channels) need to be treated separately, as depicted in Figure 5. Starting with the atmospheric flux, the neutrinos can undergo flavor change via oscillation. Since  $\nu_\tau$  production in the atmosphere is expected to be negligible at the energies relevant here, this flavor only appears through oscillation [36]. The detection rate varies between CC and NC interactions [21]. Finally, after applying the reconstruction resolutions and event classification, event counts are summed to get the final-level templates for events classified as tracks and cascades separately. Where not mentioned explicitly, the same treatment is also applied to anti-neutrinos. The final templates are the sum over both, neutrinos and anti-neutrinos.

Since the transformations computed by individual stages are independent of one another, a parameter change affecting one stage does not affect the transformations used by the other three stages, and in particular not the result of the previous stages. Therefore, we include caching functionality that reduces the overall computational expense when a number of successive templates are retrieved while changing one parameter at a time.

The transformations performed by the individual stages are dependent on the neutrino’s energy and zenith angle, and therefore must be computed and applied differentially. All stages are evaluated on a grid of points distributed over  $E_{\text{true}}$  and  $\cos\vartheta_{\text{true}}$ , with the final templates output in  $E_{\text{reco}}$ ,  $\cos\vartheta_{\text{reco}}$ , and event class. Points in energy are logarithmically spaced in the domain 1 GeV to 80 GeV while points in cosine-zenith are linearly spaced

---

<sup>11</sup>While not shown here, it is possible to extend the model with more parameters or stages to describe additional effects, such as the modeling of systematic uncertainties.

Stage	Transformation	Output
Flux	-	$400 E_{\text{true}} \times 400 \cos \vartheta_{\text{true}}$
Oscillation	$400 \times 400$	$400 E_{\text{true}} \times 400 \cos \vartheta_{\text{true}}$
Detection	$400 \times 400$	$200 E_{\text{true}} \times 200 \cos \vartheta_{\text{true}}$
Reconstruction	$200 \times 200 \times 40 \times 40 \times 2$	$40 E_{\text{reco}} \times 40 \cos \vartheta_{\text{reco}} \times 2 \text{ classes}$

Table 2: Gridpoints chosen for the staged approach in this work. The output of one stage is the input to the next stage, and the result of the detection transformation is downsampled from  $(400 \times 400)$  to  $(200 \times 200)$  by summing non-overlapping sets of  $2 \times 2$  adjacent points. Outputs of flux, oscillation, and detection are in the domain  $E_{\text{true}} \in (1, 80)$  GeV and  $\cos \vartheta_{\text{true}} \in (-1, 1)$  while the output of reconstruction is in the domain  $E_{\text{reco}} \in (1, 80)$  GeV,  $\cos \vartheta_{\text{reco}} \in (-1, 1)$ , and  $\text{class} \in \{\text{track}, \text{cascade}\}$ . Within their respective domains, points in energy are logarithmically spaced while points in cosine-zenith are linearly spaced.

between  $-1$  and  $1$ . The number of bins in each stage (for input, transformation, and output) is adjusted to reduce numerical integration errors and to avoid smearing out the physical effects under study. At the same time, this number should be kept as small as possible to reduce the computational load. An overview of the binning scheme we have employed, suitably mediating between these two effects, is given in Table 2.

The fundamental motivation for splitting up the process of template generation into a sequence of stages is that smoothing methods can be chosen for each stage that accurately reflect their unique physics, which in our example analysis apply to the detection and reconstruction stages. This approach reduces the required MC statistics with no loss of detail in the flux and flavor oscillation modeling. In contrast, smoothing events at the final level, as the traditional direct KDE does, acts on a convolution of effects, including the rapidly-varying behavior in the underlying oscillation physics. As will be shown later in this article, this difference is key to achieving higher precision with the staged approach compared to our reference methods.

For the staged approach, we emphasize that our choice of smoothing techniques is not unique. The specific techniques we employ are motivated by the typical shapes of the distributions characterized and have been found to be reliable and robust at modest computational costs. They should thus be seen as effective but non-exclusive solutions to problems of the kind discussed in this article.

Note that, in addition to the MC-based calculation of the transformations provided by the detection and reconstruction stages, we have implemented the option to produce transformations using the parametric functions of the toy model defined in Appendix A. The template produced in this way is what we refer to as “truth”.

All MC events we use with the staged approach are samples of the unbinned distributions of the toy model and are shared between the detection and reconstruction stages. For each combination of neutrino type and interaction type (for example  $\nu_e$  CC,  $\bar{\nu}_\mu$  NC), we draw an identical number of events. This number, one twelfth of the total number of events constituting a given random sample of the toy model, is referred to as the *sample size*. A given sample is used together with the event-by-event MC weighting technique to generate templates for all possible values of  $\theta$ , that is, to calculate the associated expected counts in all bins of each final-level template.



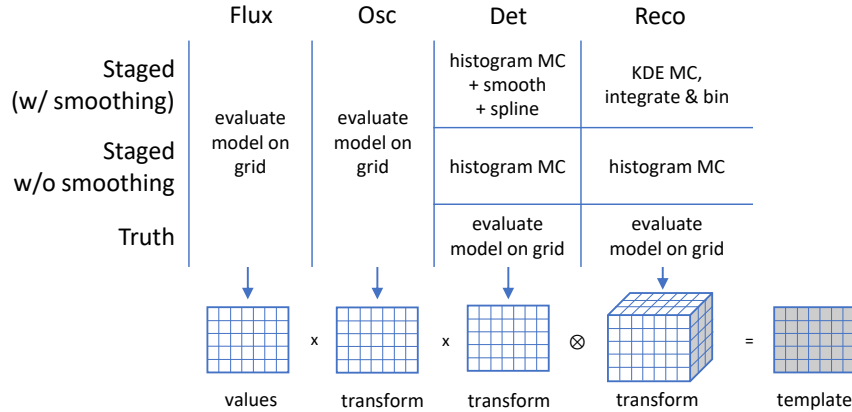


Figure 6: Operating principles of the different staged approach modes, which differ in how we generate the transformations of the last two stages. The staged approach without smoothing is employed for validation purposes in Section 6.1. See text for details.

A complete overview of the different operation modes of the staged approach is given in Figure 6, which highlights the stages at which these differ in the template generation process.

## 5. Technical Implementation of Stages

The stages within our approach, as summarized in Section 4 and illustrated in Figure 4 are subject to different technical and computational challenges due to the physics effects captured by each one. In this section we examine specific implementation details which highlight how each stage balances performance and precision requirements—even in the presence of low MC statistics.

Therefore, we include caching functionality that reduces the overall computational expense when a number of subsequent templates are retrieved while changing one parameter at a time.

### 5.1. Flux

In order to preserve the integral of a tabulated set of data, a spline is fit to the *integral* of the data rather than to the values themselves. Interpolated values in the initial space are then found by evaluating the derivative of these splines. We refer to this method as integral-preserving (IP) interpolation.

For the NMO example analysis, the tabulated data of interest are the atmospheric neutrino flux predictions from [30] provided as a function of both  $E_{\text{true}}$  and  $\cos \vartheta_{\text{true}}$ . To simplify the problem, the integration<sup>12</sup> is performed along one dimension at a time.

Consider the case with fluxes tabulated at  $M \times N$  points in  $(E_{\text{true}}, \cos \vartheta_{\text{true}})$ . To retrieve the flux at an arbitrary  $(E_{\text{true}}, \cos \vartheta_{\text{true}})$  point, say  $(x, y)$ , first one spline of the integrated

<sup>12</sup>Here, a cumulative sum of the bin values multiplied by the respective bin width.

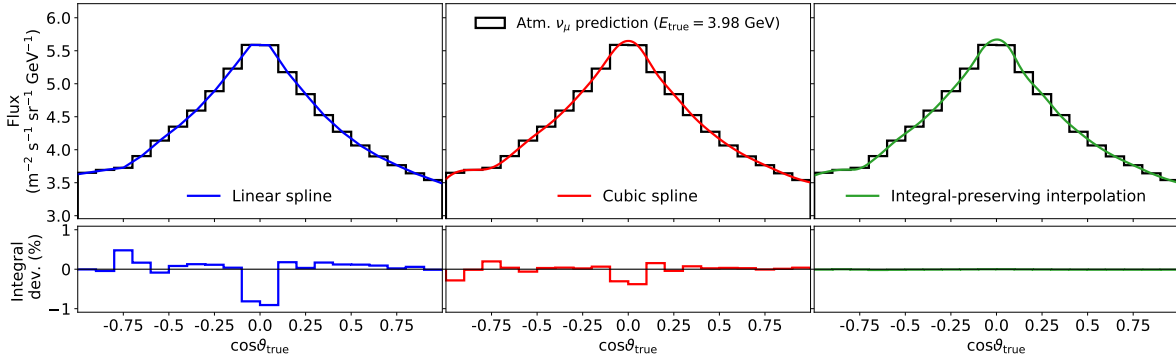


Figure 7: The top part of the figure shows three different interpolation methods applied to the same set of data points from [30] while the bottom portion shows the fractional deviation of the integral (= area under the curve) from these three methods. The deviations from the integral-preserving method presented in this paper have a maximum  $\sim 0.02\%$ .

flux as a function of  $\cos\vartheta_{\text{true}}$  is created for each of the  $M$   $E_{\text{true}}$  locations. The derivative of each of these splines is evaluated at  $y$ , yielding  $M$  flux values. The integral of these values is then interpolated with a spline, and finally this spline’s derivative is evaluated at  $x$ . This concept generalizes to higher dimensions, but can quickly become computationally intensive as the number of splines grows. While the splines used in the provided example are of one-dimensional cubic type, other spline variants or interpolation techniques can be used, as long as these allow for differentiation. For the example analysis of this article, IP interpolation is approximately an order of magnitude slower than two-dimensional cubic spline interpolation.

The IP method improves upon standard interpolation techniques in that it correctly models the turnover of the flux at the horizon ( $\cos\vartheta_{\text{true}} = 0$ ) and the behavior in the most upgoing and downgoing regions ( $\cos\vartheta_{\text{true}} \sim \pm 1$ ). This can be seen in Figure 7, which compares the results of IP to linear and cubic spline interpolation.

For the tables used in this article’s example analysis, IP interpolation preserves the integral to better than 0.5% over the complete ( $E_{\text{true}}, \cos\vartheta_{\text{true}}$ )-space. More detailed information on the IP method can be found in [37].

## 5.2. Oscillation

The oscillation library that we employ is an extension of the code described in [38], where the authors ported some of the core functions of **Prob3++** to a graphics processing unit (GPU) via the CUDA C API [39]—an application programming interface to perform general purpose computations on GPUs. We then added back in the ability to handle an arbitrary number of constant density layers of matter, allowing for highly parallel calculations of three-flavor oscillation probabilities of neutrinos that encounter a realistic radial Earth density profile, with fine-grained control over its characteristics. We implemented the oscillation calculations with floating point precision selectable to either single (32 bits, or FP32) or double (64 bits, or FP64) precision. With our code run in double precision with **Prob3++**, evaluated on a  $100 \times 100$  grid of neutrino energies  $E_{\text{true}}$  ranging from 1 GeV to 80 GeV and  $\cos\vartheta_{\text{true}}$  values

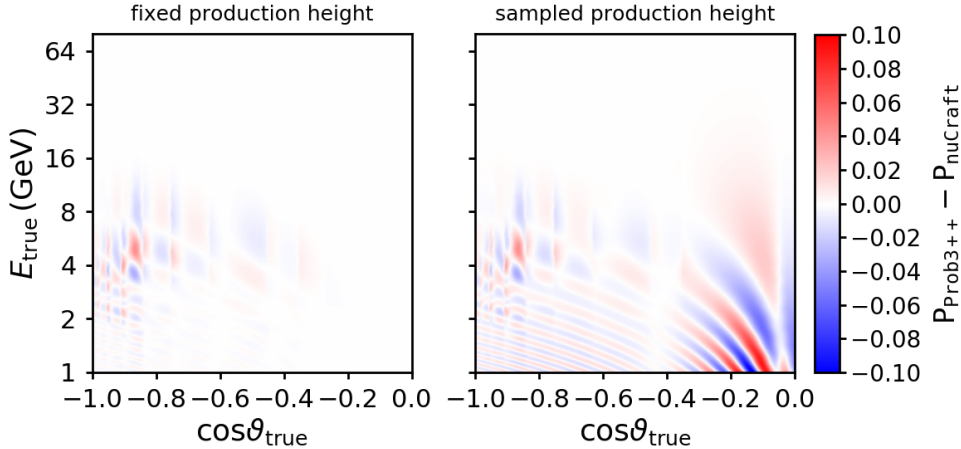


Figure 8: Deviation of  $\nu_\mu$  survival probabilities computed with **Prob3++** compared to **nuCraft**. The left panel uses a fixed production height of 20 km for both codes and twelve constant-density layers for **Prob3++**. In the right panel the values from **nuCraft** are the average probabilities for a range of neutrino production heights across the atmosphere.

spanning the region between  $-1$  and  $0$ , our GPU and CPU implementations of the **Prob3++** code produce consistent results to the level of  $10^{-14}$  or less. These differences are due to differing hardware implementations of the same mathematical operations. Switching from double to single precision on the GPU, we find that the magnitudes of the differences stay below about  $10^{-5}$  for all oscillation channels. Single precision is desirable from a performance point of view, since most GPUs comprise a larger number of single precision than double precision arithmetic units, and these extra units can be exploited by the parallelism in our code.

To evaluate the effects of an approximated Earth density profile using a limited number of constant density layers and a constant atmospheric production height—both approximations that our code makes—we compare the oscillation probabilities from our implementation of **Prob3++** against a reference model. The latter is calculated by **nuCraft** [40], which is written in Python and solves the Schrödinger equation numerically. The **nuCraft** library also supports a realistic variation of the oscillation baselines according to the distribution of atmospheric neutrino production heights described in [41] and uses an interpolated radial density profile of the Earth.

To this effect, we first fix the atmospheric neutrino production height to  $h_0 = 20$  km for both codes, and quantify the deviations arising from the coarser Earth model by calculating the  $\nu_\mu$  survival probability residuals on a fine grid in cosine zenith and energy. When approximating the Earth’s density profile with only four layers (one for each of the upper and lower mantle, and the outer and inner core), differences of up to 15% to the output of **nuCraft** are seen. These differences decrease to below 5% when using 12 density layers (see left panel of Figure 8). Using an even more detailed model with 59 layers results in differences smaller than 0.3% across the whole two-dimensional spectrum.

Comparing the 12-layer `Prob3++` probabilities to those obtained under the assumption of a more realistic distributed atmospheric production height in `nuCraft` highlights further discrepancies between the outputs of the two codes (see right panel of Figure 8). However, the largest differences ( $\sim \pm 10\%$ ) appear for near horizontal trajectories, while the residuals for  $\cos \vartheta \lesssim -0.4$  remain roughly unchanged.

Since precise modeling of both the Earth’s density profile and the atmospheric neutrino production heights come at a significant additional computational cost, depending on the analysis in question it might be desirable (and justifiable) to neglect one or both of these effects. In our example NMO analysis we find that it is sufficient to use the 12-layer model and a fixed production height. Both approximations have very little impact on the final spectra—mainly due to detector resolution effects—and since they systematically affect both NMO realizations in an almost identical manner, their effects leave the measurement comparing the two mass orderings largely unaffected. Moreover, while the atmospheric flux peaks in horizontal direction (seen, for example, in Figure 9), negligible matter effects for the corresponding trajectories result in very little intrinsic sensitivity of this part of the spectrum to the NMO.

### 5.3. Detection

As a reminder, the effective areas are quantities used to translate an incoming flux to the event rates in the detector. These effective areas are calculated from a limited number of MC events, hence they can suffer from statistical fluctuations which can be a non-negligible contribution to the total uncertainty of the final physics result. At the same time, effective areas are typically well-behaved quantities in energy and zenith angle (under some realistic assumptions, e.g., that no discontinuous selection cuts are applied and no gaps exist in the detector acceptance). Therefore, smoothing techniques can be applied to alleviate the unwanted uncertainty contributions from statistical fluctuations.

For charged current interactions, we compute the effective area separately for each neutrino flavor. In contrast, we do not distinguish between flavors for neutral current (NC) interactions, since their cross sections are identical. Neutrinos and anti-neutrinos are handled independently, accounting for a total of eight independent effective area functions. For convenience we include the multiplication by detector exposure time ( $t_{\text{exp}}$ ) in the same step, which means that this stage outputs event counts ( $N_{\text{events}}$ ) instead of rates

$$N_{\text{events}} = \Phi[\text{m}^{-2}\text{s}^{-1}] \cdot A_{\text{eff}}[\text{m}^2] \cdot t_{\text{exp}}[\text{s}], \quad (4)$$

for some input flux ( $\Phi$ ).

In our staged approach we first evaluate the effective areas on a fine grid in  $(E_{\text{true}}, \cos \vartheta_{\text{true}})$  using the MC events via MC integration, where, when generating events, the sampling is chosen to provide a relatively uniform coverage across all grid points. For our example case study, we use a uniform sampling across  $\cos \vartheta_{\text{true}}$  and a power law spectrum for the energies  $\propto E_{\text{true}}^{-1}$  to closely follow actual IceCube oscillation analyses. (Note that an optimization of the sampling choices would benefit both the staged approach and the reference methods.) Still, for small sample sizes, some grid points may have no associated events, leading to gaps

in the distribution. We remove these by applying a simple Gaussian smearing along the two-dimensional grid. In a second step, cubic splines are employed to perform smoothing. Here, first, splines are created along the  $E_{\text{true}}$  dimension individually for every  $\cos \vartheta_{\text{true}}$  bin, and evaluated to obtain new values for every grid point. Then, this splining procedure is repeated along the  $\cos \vartheta_{\text{true}}$  dimension.

Figure 9 shows the truth template of  $\nu_\mu$  CC events on a grid with  $n_{\text{bins}} = 40 \times 40$  points together with the fractional deviations that arise when the same template is obtained from MC samples<sup>13</sup> of different sizes using direct histogramming versus the smoothing method described above. We use  $\nu_\mu$  CC events as an example here and below. Table 3 gives the average (binwise, i.e. per degree of freedom)  $\chi^2$  values defined as

$$\langle \chi^2 \rangle = \frac{1}{n_{\text{bins}}} \chi^2 = \frac{1}{n_{\text{bins}}} \sum_{i=1}^{n_{\text{bins}}} \frac{(\mu'_i - \mu_i^{\text{ref}})^2}{\mu_i^{\text{ref}}} \quad (5)$$

and maximal  $\chi^2$  values defined as

$$\chi_{\text{max}}^2 = \max_{1 \leq i \leq n_{\text{bins}}} \left[ \frac{(\mu'_i - \mu_i^{\text{ref}})^2}{\mu_i^{\text{ref}}} \right] \quad (6)$$

by which the templates from our method and from direct histogramming deviate from truth (with bin counts  $\mu_i^{\text{ref}}$ ).

The  $\chi^2$  values provide direct insight into how the accuracy of the template description compares to the statistical uncertainty of the real data that would be observed. Since the observed data underlies Poisson fluctuations it has an average deviation from truth of  $\chi^2 = 1$  per bin. An analysis of real data, however, can only test templates based on MC. These exhibit their own statistical uncertainties, resulting in finite  $\chi^2$  deviations from truth, shown in Table 3 as a function of MC sample size. It is essential that these inaccuracies inherent to the template generation process are considerably smaller than the statistical fluctuations in data in order to ensure accurate statistical inference.

Applying our method we find deviations that are lower by a factor of about 40 for the smallest MC set, and by a factor of about 13 for the largest. It is noteworthy that the maximum deviation ( $\chi_{\text{max}}^2$ ) across all bins decreases monotonically with MC sample size, confirming that the used smoothing method does not introduce any observable bias.

#### 5.4. Reconstruction

The usual way to obtain templates in the space of reconstructed variables is to place each individual MC event in the final-level distributions according to the reconstruction information that the event carries. This is the case for both methods that are used for comparison: direct histogramming and direct KDE, the only difference between these being how the final-level distributions are estimated. While this approach correctly takes into account joint dependencies of the event reconstruction on the involved variables, it is particularly

---

<sup>13</sup>Generated from the toy model in Appendix A.

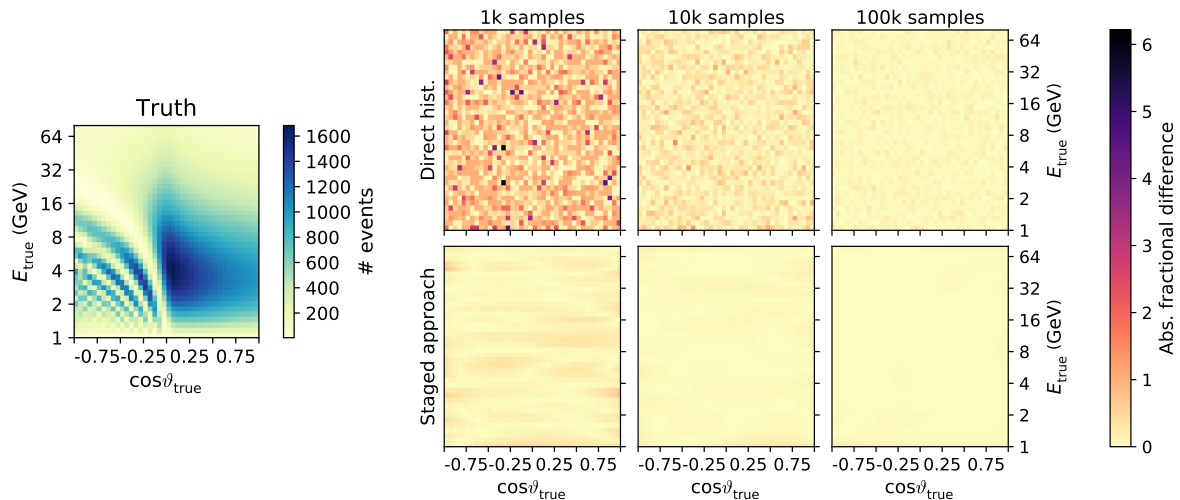


Figure 9: Parametric reference distribution after the first three stages (flux, oscillation, and detection) for the  $\nu_\mu$  CC channel in  $(\cos \vartheta_{\text{true}}, E_{\text{true}})$  (left panel) and relative residuals  $(|N - N_{\text{true}}|/N_{\text{true}})$  for the direct histogramming (right panel, top row) and our proposed method (right panel, bottom row) on a  $40 \times 40$  grid using different amounts of simulated events. (Note that the numbers are not percentages.) The three columns in the right panel represent different MC sample sizes of  $10^3$ ,  $10^4$ , and  $10^5$  events, respectively. The samples are drawn from the unbinned toy model distributions of Appendix A.

sensitive to small MC sample sizes due to the potentially high dimensionality of the space of reconstructed variables. In contrast, the staged approach uses the available MC simulation to construct detector resolution functions which we integrate to form a transformation that maps a template in true variables (such as that shown on the left in Figure 9) onto the space of reconstructed variables, what we refer to as the final-level template.

In the case study of the NMO analysis, the mapping of true variables ( $E_{\text{true}}$  and  $\cos \vartheta_{\text{true}}$ ) to reconstructed variables ( $E_{\text{reco}}$ ,  $\cos \vartheta_{\text{reco}}$ , and event class) is extracted from the MC as a “migration” tensor of order five,  $\mathcal{M}_{ijklm}$ . It maps the histogram of event counts in the two-dimensional space of true variables,  $h_{ij}$ , to the observed histogram of event counts in the

Sample size		$10^3$	$10^4$	$10^5$	$10^6$
Direct hist.	$\langle \chi^2 \rangle$	215	22.5	2.07	0.201
	$\chi^2_{\text{max}}$	21600	1810	79.4	11.2
Staged approach	$\langle \chi^2 \rangle$	5.14	0.526	0.0615	0.0156
	$\chi^2_{\text{max}}$	460	17.2	2.27	0.975

Table 3: Average  $\chi^2$  per bin and the worst-case bin’s  $\chi^2$  value comparing templates on a  $40 \times 40$  grid in  $(E_{\text{true}}, \cos \vartheta_{\text{true}})$ -space (i.e., before applying reconstruction resolutions) generated by direct histogramming (top) and the smoothed-staged approach (bottom) with the toy model’s reference template. Shown are values obtained for independent input MC samples of various sizes (from  $10^3$  up to  $10^6$  events per flavor/interaction type).

three-dimensional space of reconstructed variables,  $h'_{klm}$ :

$$h'_{klm} = \sum_{i,j} \mathcal{M}_{ijklm} h_{ij} . \quad (7)$$

The reconstruction transform in general has to be computed as a five-dimensional transform, as all five dimensions can depend on one another—i.e. they are correlated. Studying the correlations among the dimensions in our particular MC revealed, however, that  $E_{\text{reco}}$  only depends on event class and  $E_{\text{true}}$ ,  $\cos \vartheta_{\text{reco}}$  depends on event class and both input dimensions, and event class only depends on  $E_{\text{true}}$ . For each of the three reconstruction variables, we subdivide the MC in the quantity’s dependent dimensions to the point that correlations are not visible and that all events in the subdivision can be assumed to be samples from the same one-dimensional distribution.—i.e. the resolution functions we generate.

There is a trade-off in terms of how much to subdivide the MC for producing these resolution functions. Since resolution changes as a function of a dependent dimension, sufficiently narrow subdivisions in that dimension group together MC events drawn from essentially the same distribution. Subdivisions that are too wide will group together events drawn from different distributions and the resulting resolution functions will be erroneous. However, narrower subdivisions admit fewer MC events in each subdivision and so lead to greater statistical variations in the estimated resolution functions (i.e., their shapes will be more affected by random fluctuations in the MC).

To balance these competing factors, we devised the following heuristic. For the quantity being characterized, we divide each dependent dimension evenly—except event class, which is binary.  $E_{\text{true}}$  is divided evenly in log-space to help ensure even subdivisions group together events with similar energy resolution, as this quantity changes more rapidly at low  $E_{\text{true}}$  than at high  $E_{\text{true}}$ . We allow each subdivision of  $E_{\text{true}}$  to separately expand enough to capture at least 100 events, and at least 500 events in each subdivision of  $\cos \vartheta_{\text{true}}$ . If expansion is performed, subdivisions’ upper and lower edges are expanded by the same factor (up to the limits of the dimension). The captured events are then used to produce resolution functions.

The remaining parameters that require tuning in this heuristic are the number of subdivisions to use for each dependent dimension for each quantity being characterized. For this, we visually inspect the 2-dimensional distributions of each characterized quantity as a function of each dependent dimension and require that the events in each subdivision do not display strong dependence on the dependent dimension.

If the functional form of the resolution functions is known, a parametric model of this form fit to the MC yields the most accurate and lowest variance reconstruction transform. However, as we do not know the form of these functions, a non-parametric density estimation technique is used to approximate them. In particular, we chose to use adaptive KDE [42] with bandwidths scaled uniformly such that the narrowest is that found from the Improved Sheather Jones (ISJ) algorithm [43]. KDE works by placing a kernel function (we use a Gaussian) centered at the value of each event’s variable to be described and then summing over all kernels. Adaptive bandwidth KDE uses different widths for each kernel, where the bandwidths are inversely proportional to the density of points near the location of the kernel. The ISJ bandwidth selection algorithm used to normalize the kernel widths is an improvement

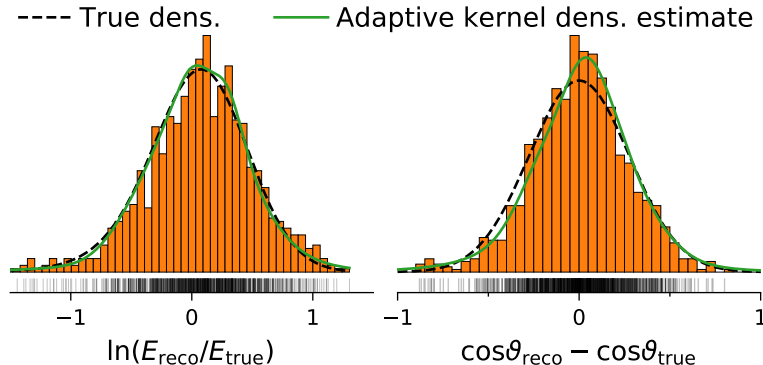


Figure 10: Example energy and cosine-zenith-angle resolution distributions for  $\nu_\mu$  CC events classified as cascades, estimated with histograms and adaptive KDE. Energy resolution is shown for 100 events with  $E_{\text{true}} \in [26.7, 29.8]$  GeV and cosine-zenith resolution for 100 events with  $E_{\text{true}} \in [1.0, 1.1]$  GeV. The samples used to construct the histogram and KDE are shown by vertical lines beneath the histograms.

over predecessor algorithms (e.g., [44, 13]) in that it does not make assumptions that the quantity being estimated is drawn from a Gaussian distribution. In our experience, this outperforms fixed bandwidth KDE by not underestimating the heavy-tailed distributions we encounter, but it bears repeating that other density estimation techniques can yield better or worse results depending on the specifics of the MC in question. An example of two resolution functions (one for both energy and zenith angle, respectively) estimated using the adaptive KDE method is shown in Figure 10.

Figure 11 again demonstrates that templates obtained from our KDE-based reconstruction stage deviate much less from the parametric reference template after reconstruction than templates from direct histogramming of reconstructed MC events.

## 6. Validation and Comparison of Templates

This section more closely examines the templates generated with the staged approach and compares them—along with those generated by the other two methods (histograms and KDE)—to the parametric reference distributions of the toy detector model. This validation is split into two parts. The first examines the grid of points that are used to numerically approximate the integral over the first three stages, whereas the effect of smoothing is investigated in the second.

### 6.1. Sampling Grid

In order to demonstrate the validity of our choice of grid points shown in Table 2 as well as the equivalence between the staged approach and traditional MC weighting as grid point spacing in  $E_{\text{true}}$  and  $\cos\vartheta_{\text{true}}$  is reduced, we compare the staged approach without smoothing (i.e. using raw histograms as transforms in place of smoothing functions and KDEs) to direct histogramming. The specific comparison done here without smoothing is solely for the purpose of validating the principle of stages vs. direct histograms.



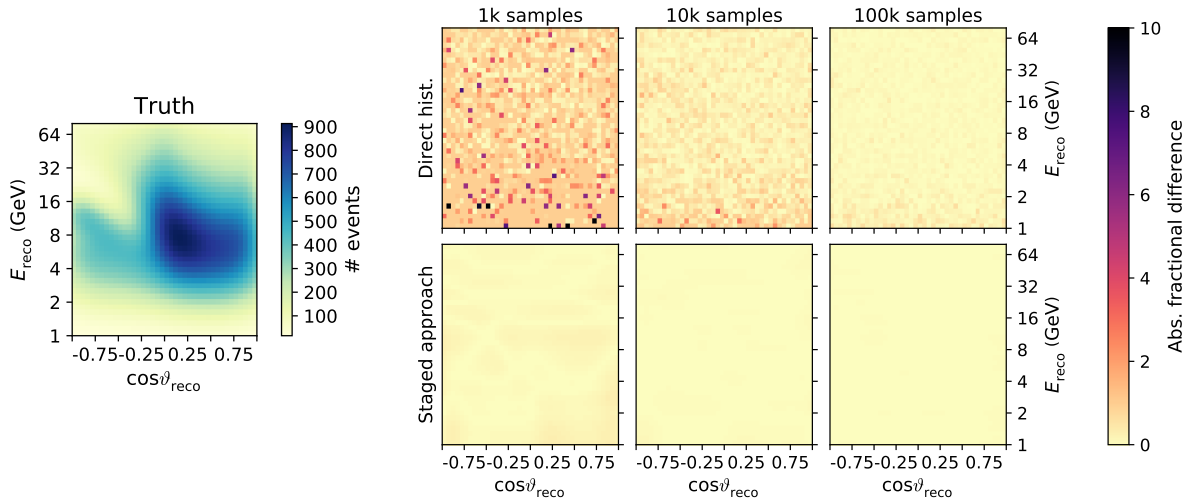


Figure 11: Same as Figure 9, but comparing final-level templates after all four stages are applied. Note that the residuals in the 1k-samples plot for direct histogramming go up to 31 but the scale is clipped at 10.

Grid ( $M \times N$ )	$40 \times 40$	$80 \times 80$	$160 \times 160$	$320 \times 320$	$640 \times 640$	$1280 \times 1280$
$\langle \chi^2 \rangle$	0.01067	0.00253	0.00060	0.00014	0.00003	0.00001
$\chi^2_{\max}$	1.45906	0.46930	0.19718	0.04974	0.00634	0.00172

Table 4: Average and maximal  $\chi^2$  deviations per bin of the final  $40 \times 40 \times 2$  binning between final templates of non-smoothed staged approach and direct histogramming, for different grid point densities in  $(E_{\text{true}}, \cos \vartheta_{\text{true}})$  for the first three stages, using an MC sample of  $10^6$  events. The last (=reconstruction) stage uses a reduced binning, as described in the text.

Table 4 shows the  $\chi^2$  difference (cf. Equations (5) and (6)) between the final templates obtained from the staged approach (with bin counts  $\mu'_i$ ) and direct histogramming ( $\mu_i^{\text{ref}}$ ) for a variety of grid point densities in  $E_{\text{true}}$  and  $\cos \vartheta_{\text{true}}$ , using the same MC set of size  $10^6$  for both methods. These templates are output with a binning of  $40 \times 40 \times 2$  in  $E_{\text{reco}}, \cos \vartheta_{\text{reco}}$ , and event class. The relative decrease in the average  $\chi^2$  value roughly scales with the inverse of the relative grid density increase, thus confirming that the two methods will agree to arbitrary precision in the asymptotic limit. In the following, for practical reasons we limit ourselves to the specific case summarized in Table 2.

## 6.2. Smoothing

To validate the final templates with smoothing applied at each stage, we compare them directly to truth. For reference, we also show the agreement resulting from both the direct histogramming and the direct KDE methods.

While Table 5 quantifies deviations from the reference distributions again in terms of  $\chi^2$  and in dependence of MC sample size, Figure 12 displays the final-level templates for each of the aforementioned methods using a sample with  $10^4$  events.

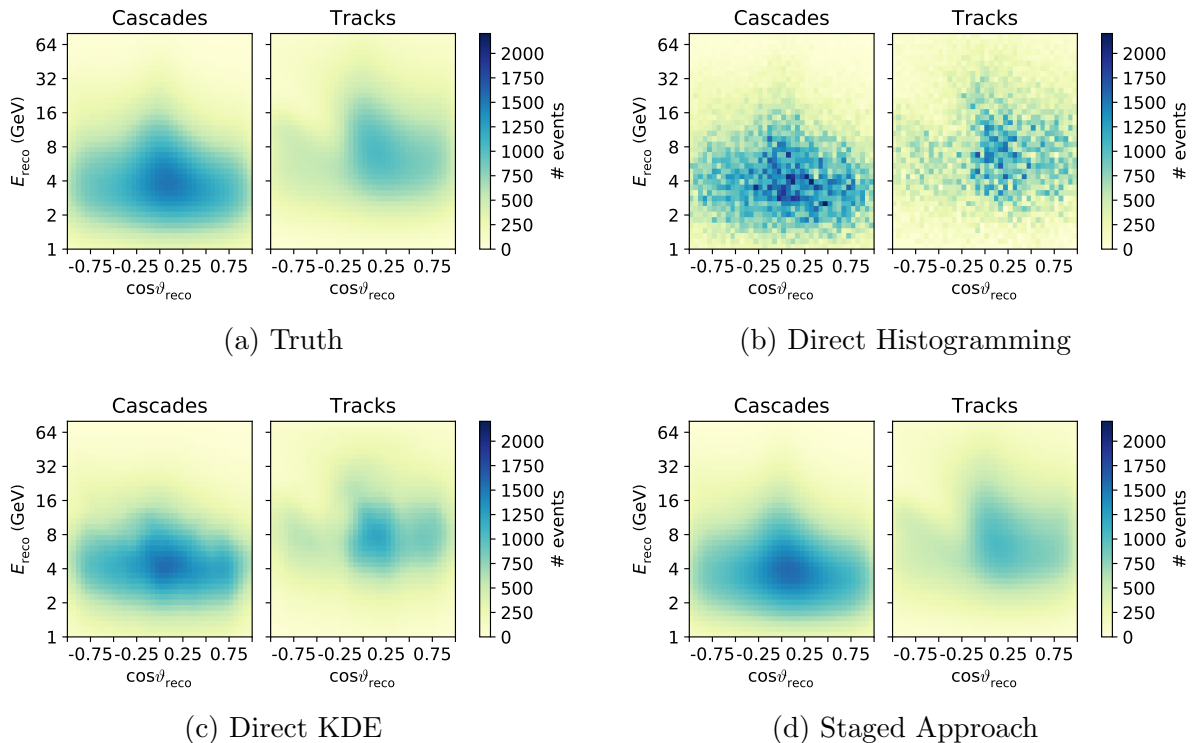


Figure 12: Final-level templates used for the example data analysis. The reference distributions (truth) obtained directly for the toy detector model parameterizations are shown in panel (a). Given the same sample of  $10^4$  events the estimated distributions using histograms are shown in panel (b), using KDEs in panel (c), and using the staged approach in panel (d).

The staged approach outperforms the two alternatives in terms of  $\chi^2$  values by more than one order of magnitude for all those sample sizes studied here. Furthermore, inaccuracies of the templates from the staged approach scale with the inverse of sample size almost as fast as those of templates from direct histogramming. In addition, it is noteworthy that the KDE method shows comparably slow convergence, i.e., it performs worse than direct histogramming for the sample size of  $10^6$ .

While for the experimental data (or pseudo-data) one expects statistical fluctuations on the order of  $\chi^2 = 1.0$  per bin, the accuracy of the templates must be better than this. As shown in Table 5, considering a sample size of  $10^4$  and the staged approach, the average  $\chi^2$  deviation from truth (using the same  $\chi^2$  definition as for data) is only about 30% of what is expected just from statistical fluctuations in data, while more than  $10^6$  events would be necessary to achieve the same average  $\chi^2$  using direct histogramming or KDE. (See Table 5 for details.) Therefore, to reach an equal accuracy, two or more orders of magnitude larger samples are needed for histogramming or KDE compared to the staged approach. The next section illustrates the implications for running a data analysis.

Sample size		$10^3$	$10^4$	$10^5$	$10^6$
Direct Histogramming	$\langle\chi^2\rangle$	468	42.6	4.27	0.458
	$\chi_{\max}^2$	$3.4 \cdot 10^4$	906	138	10.5
Direct KDE	$\langle\chi^2\rangle$	32.2	11.4	3.67	1.25
	$\chi_{\max}^2$	245	90.2	50.3	25.3
Staged Approach	$\langle\chi^2\rangle$	3.01	0.303	0.111	0.0301
	$\chi_{\max}^2$	47.4	3.03	1.80	0.387

Table 5: Average and maximal  $\chi^2$  deviations per bin of the final  $40 \times 40 \times 2$  binning between final templates of the three shown methods and truth, for independent input MC samples of various sizes. Note that the staged approach has smoothing applied (the default), in contrast to Table 4.

## 7. Example Analysis Results

To illustrate the impact of sample size, we show the resulting  $\sqrt{\Delta\chi^2}$  as an estimate for the sensitivity to the NMO for our example analysis in Figure 13. For reference, the true result is derived directly from the exact templates based on the parametric toy detector model and lies at  $\sqrt{\Delta\chi^2} = 5.75$ . For the three methods discussed throughout this paper, the statistical uncertainty of the obtained sensitivity is indicated by error bars in the figure. This uncertainty is computed from several statistically independent MC sets<sup>14</sup> and indicates the central 68% quantile of each ensemble. In particular, as the sensitivity proxy does not take into account MC uncertainty [5, 6], this range is not, a priori, expected to reflect any sensitivity bias for the three methods.

The uncertainty reveals that the methods exhibit quite different intrinsic fluctuation of their respective sensitivity estimates, as well as different scaling behavior of the variance with sample size. As sample size decreases, direct histogramming without any smoothing applied results in an increasing overestimation of a VLV $\nu$ T’s ability to exclude the wrong neutrino mass ordering. In the most extreme case shown here (corresponding to the smallest sample size of  $10^3$ ), the sensitivity is estimated to be more than one order of magnitude greater than the actual capability of the experiment. Only for the sample size of  $10^7$  does direct histogramming indeed give reliable results. This is expected from the simple rule of thumb (cf. Section 2.2) of  $\mathcal{O}(10^4)$  events per bin  $\times \mathcal{O}(10^3)$  bins.

Illustratively, an undersampling of the detector response distributions due to low MC statistics is highly likely to lead to an overestimation of the experiment’s sensitivity because the NMO signature that is present in the space of true variables is carried over to random bins in the reconstructed observables with reduced cancellation<sup>15</sup>.

Applying KDE smoothing to the weighted events instead of histogramming them (i.e., direct KDE) leads to a reduction of the overestimated sensitivity for sample sizes of up to at

<sup>14</sup>Each MC set is used together with the staged approach to generate one Asimov toy data template and  $\mathcal{O}(10^3)$  “test” templates.

<sup>15</sup>For example, if a bin in the final-level template is solely populated by (unweighted) MC neutrinos, and no anti-neutrinos, or vice-versa, it will contribute artificially strong to the overall NMO sensitivity due to the missing summation over both event types (cf. Section 3).

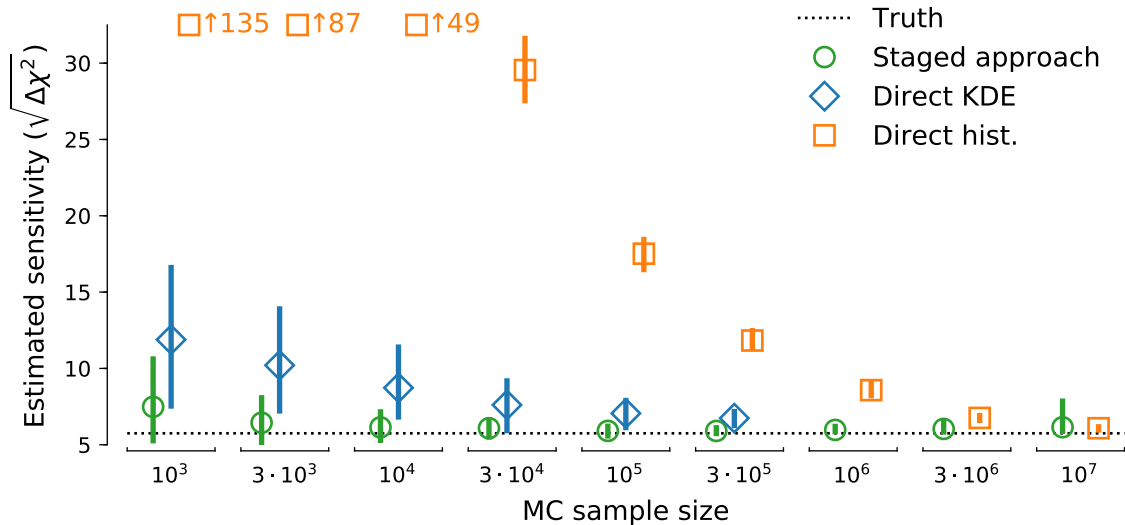


Figure 13: Estimated sensitivity ( $\sqrt{\Delta\chi^2}$ ) to the NMO vs. sample size for direct histogramming, direct KDE, and the proposed staged smoothing methods applied to multiple (between 50 and 200) statistically independent toy MC sets. Vertical lines indicate central 68% quantiles of the ensemble. The dashed horizontal line shows the significance obtained from truth templates based on the parametric toy detector model. The staged approach outperforms the other methods—both in terms of bias and variance—for sample sizes through  $3 \cdot 10^6$ , with direct histogramming only matching the staged approach using  $10^7$  samples. Note that no data points exist for direct KDE and sample sizes above  $3 \cdot 10^5$ , as computational processing times become impractically large. Also note that direct histogramming is off-scale high for fewer than  $3 \cdot 10^4$  events (mean values are indicated to the right of the corresponding markers).

least  $3 \cdot 10^5$  but is not able to eliminate the bias entirely for the tested sample sizes. For sample sizes larger than  $\mathcal{O}(10^5)$ , the direct KDE method is too computationally expensive to deliver results within a reasonable time (for more details on timing, see Section 8).

The estimated sensitivity using the staged approach is statistically compatible with the true sensitivity across the whole range of sample sizes considered. It shows no bias and lower variance for predicting sensitivity to physics compared to the other methods within the limits of our testing.

## 8. Benchmarks

Whether a given analysis method is useful in a realistic setting depends not only on its numerical reliability, but also on how long it takes to compute the quantity of interest (note that this duration is in addition to the initial time needed to generate the MC). For reference, we performed benchmarks of the template generation times in the course of a typical analysis process<sup>16</sup>. These are compiled in Figure 14.

<sup>16</sup>Timings were obtained on a computer with an Intel Xeon E5-1660 v3 3.0 GHz CPU and an NVIDIA GeForce GTX Titan X GPU.

Note that no initial start-up times—such as the construction of the smearing kernels used within the reconstruction stage—are included here. For all three methods separate timings based on our CPU-only and GPU-accelerated implementations are provided.

While for sample sizes below  $10^4$  to  $10^5$  events direct histogramming is the fastest method, it is unusable owing to the large fluctuations associated with the templates it produces, which in turn result in the grossly overestimated sensitivities shown in Figure 13. Direct KDE only proves competitive when used in conjunction with the smallest datasets. The faster-than-linear scaling of its computational needs with sample size then quickly renders it impractical to use. Our proposed method is independent of sample size by construction (excluding initial start-up costs), but will get more expensive if a finer grid point spacing is desired.

The timing difference between the CPU and GPU implementation of the staged approach is not as large as for the other methods, since it is only using the GPU for parallelization of the neutrino oscillation weights calculation (whereas the other methods make use of the GPU more extensively).

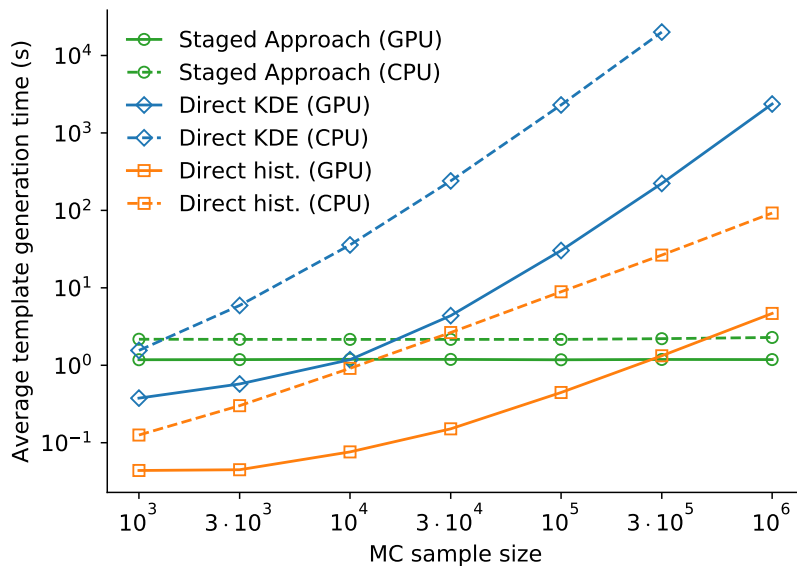


Figure 14: Average template generation time during a typical analysis for input datasets of varying size, shown for the direct histogramming, the direct KDE, and the staged approach. Solid lines represent timings based on (partial) GPU acceleration, whereas the dashed ones are for CPU-only calculations.

## 9. Summary

The search for small physics effects in high statistics neutrino oscillation experiments requires careful treatment and use of simulated data. Statistical fluctuations within distributions obtained from Monte Carlo simulations are able to severely distort an analysis, rendering derived constraints or sensitivities essentially meaningless.

The staged approach we have presented serves two main purposes. Firstly, computational expense is reduced through sampling of physics and detector response distributions on a discrete grid instead of computing a weight for every individual Monte Carlo event. In this respect, we have demonstrated that our method of breaking down the template generation into independent stages converges to the MC weighting scheme when using a grid of a high enough, albeit feasible, density. For a fixed number of grid points, the template generation time has been shown to be independent of the input sample size. Secondly, the staged approach allows the application of smoothing techniques to a detector’s response functions. In the specific example shown, the detection stage characterizes the detector’s effective area by integrating weighted MC events on a grid and smoothing the resulting histogram, followed by the event reconstruction stage using an adaptive KDE smoothing on the resolution functions applied to arrive at final-level templates. This has proven superior to the smoothing of the final event distributions since it is faster and—even more importantly—yields more accurate and robust results. The presented choice of smoothing techniques works sufficiently well for our purposes, but this choice is neither unique nor do we claim it to be optimal, and it depends on the wider experimental context. Beside this choice, our overall approach may prove particularly useful when a fast assessment of the physics potential of various detector designs is desired, or when analysis methodologies are optimized. Any final-level analysis will likely rely on large quantities of MC to guarantee the precise and accurate modelling of the experiment.

In the example neutrino mass ordering analysis that we have conducted—to benchmark and compare the different approaches—we found that direct histogramming of events leads to a gross overestimation of sensitivities when used in conjunction with small numbers of events ( $\lesssim 10^6$  events for our toy model). Conversely, the proposed staged approach leads to correct results that are largely unaffected by the sample size across the tested range and the variance of results is small compared to the result above about  $10^4$  neutrino events. This means that the necessary amount of simulated events is reduced significantly (by about two orders of magnitude in our example)—an important aspect especially since Monte Carlo event simulation and reconstruction times can represent major hurdles to progress in the field of neutrino oscillation experiments.

## Acknowledgments

The IceCube collaboration gratefully acknowledges the support from the following agencies and institutions: USA – U.S. National Science Foundation-Office of Polar Programs, U.S. National Science Foundation-Physics Division, Wisconsin Alumni Research Foundation, Center for High Throughput Computing (CHTC) at the University of Wisconsin-Madison, Open Science Grid (OSG), Extreme Science and Engineering Discovery Environment (XSEDE), U.S. Department of Energy-National Energy Research Scientific Computing Center, Particle astrophysics research computing center at the University of Maryland, Institute for Cyber-Enabled Research at Michigan State University, and Astroparticle physics computational facility at Marquette University; Belgium – Funds for Scientific Research (FRS-FNRS and FWO), FWO Odysseus and Big Science programmes, and Belgian Federal Science Policy Of-

fice (Belspo); Germany – Bundesministerium für Bildung und Forschung (BMBF), Deutsche Forschungsgemeinschaft (DFG), Helmholtz Alliance for Astroparticle Physics (HAP), Initiative and Networking Fund of the Helmholtz Association, Deutsches Elektronen Synchrotron (DESY), and High Performance Computing cluster of the RWTH Aachen; Sweden – Swedish Research Council, Swedish Polar Research Secretariat, Swedish National Infrastructure for Computing (SNIC), and Knut and Alice Wallenberg Foundation; Australia – Australian Research Council; Canada – Natural Sciences and Engineering Research Council of Canada, Calcul Québec, Compute Ontario, Canada Foundation for Innovation, WestGrid, and Compute Canada; Denmark – Villum Fonden, Danish National Research Foundation (DNRF), Carlsberg Foundation; New Zealand – Marsden Fund; Japan – Japan Society for Promotion of Science (JSPS) and Institute for Global Prominent Research (IGPR) of Chiba University; Korea – National Research Foundation of Korea (NRF); Switzerland – Swiss National Science Foundation (SNSF); United Kingdom – Department of Physics, University of Oxford.

## Appendix A. Toy Data Model

In the following we provide a parametric toy detector model used to transform the oscillated atmospheric fluxes into event counts. The functions we use either serve as direct inputs (truth) to the various stages of the simulation chain laid out in Section 4, or as sampling distributions from which toy MC samples are drawn. We point out here that these are entirely empirically motivated, and should only be seen as proxies of the performance indicators in next-generation VLV $\nu$ Ts (such as the IceCube Upgrade [1], PINGU [2, 3], or KM3NeT/ORCA [4]).

Simplifications or limitations of the model do not affect the computational analysis techniques themselves. Rather, the goal in the following is to capture the most essential features of the expected detector response: threshold effects in detection, the finite accuracy and skew of reconstruction resolution functions, as well as limited flavor and charge identification capabilities. This does not invalidate the conclusions drawn from comparing the various analysis approaches.

### Appendix A.1. Detection Efficiency

We assume a detector of fiducial mass  $M_{\text{fid}} = 10$  megaton, with a neutrino detection energy threshold of  $E_{\text{th}} = 1$  GeV for all neutrino flavors and interaction channels apart from  $\nu_\tau$  charged current (CC) interactions, where the intrinsic interaction threshold is higher, at  $E_{\text{th}} = 3.5$  GeV. The detector’s effective mass  $M_{\text{eff}}^\alpha = \rho_{\text{ice}} V_{\text{eff}}^\alpha$  for a given combination,  $\alpha$ , of flavor and interaction type, where  $\rho_{\text{ice}}$  is the ice density and  $V_{\text{eff}}^\alpha$  the detector’s corresponding effective volume, exhibits a phenomenological dependence on true neutrino energy,  $E_{\text{true}}$ , asymptotically approaching  $M_{\text{fid}}$  according to an exponential function:

$$M_{\text{eff}}^\alpha(E_{\text{true}}) = M_{\text{fid}} \times \left(1 - e^{-k_\alpha \times (E_{\text{true}}/\text{GeV} - E_{\text{th}}/\text{GeV})}\right) \text{ for } E_{\text{true}} > E_{\text{th}} . \quad (\text{A.1})$$

We include three effective masses to cover all the neutrino interaction channels: one for  $\nu_e$ ,  $\bar{\nu}_e$ ,  $\nu_\mu$ , and  $\bar{\nu}_\mu$  CC, one for  $\nu_\tau$  and  $\bar{\nu}_\tau$  CC, and one for all NC channels. For the CC

channels we choose  $k_\alpha = 0.4$ , while for the NC channels the function rises more slowly, with  $k_\alpha = 0.1$ . The left panel of Figure A.15 shows these dependencies for neutrino energies up to  $E_{\text{true}} = 80$  GeV. The detector can be roughly considered fully efficient ( $M_{\text{eff}} = M_{\text{fid}}$ ) for all detection channels above  $E_{\text{true}} \approx 50$  GeV.

The  $\nu$ - $\bar{\nu}$  asymmetry—which is required to make the NMO measurement—will be introduced through differences in flux and cross sections, i.e., it will become apparent in the detector’s effective area. The latter we obtain from the effective mass via the conversion

$$A_{\text{eff}}^\alpha(E_{\text{true}}) = \sigma_\alpha(E_{\text{true}}) \times n_{\text{ice}}/\rho_{\text{ice}} \times M_{\text{eff}}^\alpha(E_{\text{true}}), \quad (\text{A.2})$$

where  $\sigma_\alpha$  is the total neutrino-nucleon cross section for a given flavor-interaction channel  $\alpha$ ,  $n_{\text{ice}} \approx 6 \times 10^{23} \text{ cm}^{-3}$  is the nucleon density in ice, and  $\rho_{\text{ice}} \approx 0.92 \text{ g cm}^{-3}$  the mass density.

We also make some simplifying assumptions about the cross sections used in Equation (A.2), in that we take  $\nu_e$  and  $\nu_\mu$  ( $\bar{\nu}_e$  and  $\bar{\nu}_\mu$ ) CC cross sections to be the same, as well as all  $\nu_x$  ( $\bar{\nu}_x$ ) NC cross sections. In addition, we model all the mentioned cross sections to rise strictly linearly with  $E_{\text{true}}$  above  $E_{\text{true}} = 1$  GeV [45]:

$$\sigma_\alpha(E_{\text{true}})/E_{\text{true}} = c_\alpha \times 10^{-38} \text{ cm}^2 \text{ GeV}^{-1}, \quad (\text{A.3})$$

where we set

$$c_{\nu_{e,\mu} \text{ CC}} = 2c_{\bar{\nu}_{e,\mu} \text{ CC}} = 0.70, \quad (\text{A.4})$$

$$c_{\nu_x \text{ NC}} = 2c_{\bar{\nu}_x \text{ NC}} = 0.25. \quad (\text{A.5})$$

To obtain  $\nu_\tau$  ( $\bar{\nu}_\tau$ ) CC effective areas, we interpolate the corresponding neutrino-nucleon cross section curves given in [46]. All resulting effective areas as a function of neutrino energy are depicted in the right panel of Figure A.15. We take these to be invariant in azimuth, but universally introduce an arbitrary, smooth polynomial modification  $M$  with the zenith angle dependency

$$M(x) = \frac{1}{20}(-x^3 + x^2 - x) + 1 \quad (x \equiv \cos \vartheta_{\text{true}}), \quad (\text{A.6})$$

which we normalize to unit area<sup>17</sup>.

### Appendix A.2. Reconstruction Resolutions

Neutrino zenith resolutions with respect to  $\cos \vartheta$  are represented by single Gaussian distributions. The distributions’ parameters are taken as functions of  $E_{\text{true}}$  only. For each flavor and interaction channel, we assign a mean  $\mu_{\Delta \cos \vartheta}(E_{\text{true}}) = 0$  across all energies, and a standard deviation of  $\sigma_{\Delta \cos \vartheta}(E_{\text{true}}) = \frac{0.3}{\sqrt{E_{\text{true}}/\text{GeV}}} + 0.05$ .

Neutrino energy resolutions we describe using right-skewed Gumbel distributions. These are shifted and scaled by  $\mu'$  and  $\sigma'$  with respect to their standard form, via the transformation

---

<sup>17</sup>  $A_{\text{eff}}(E_{\text{true}})$  is the average over the full sky,  $\cos \vartheta_{\text{true}} \in [-1, +1]$ .



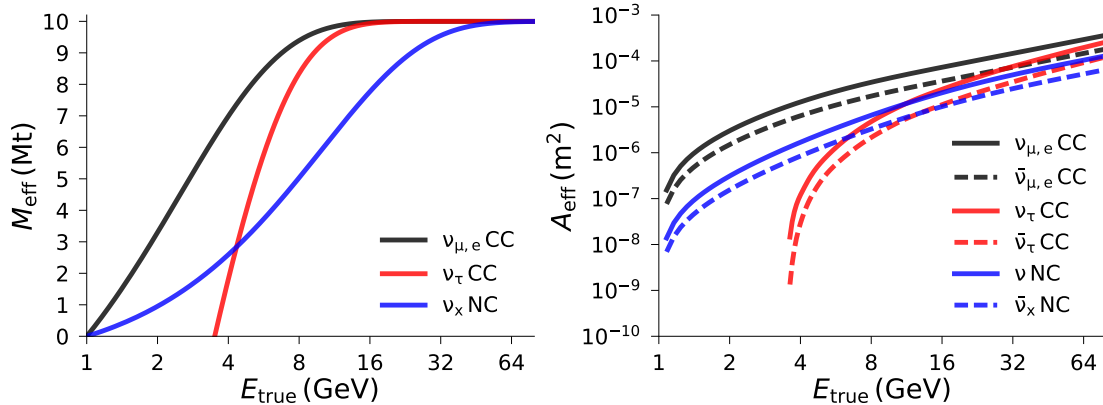


Figure A.15: Effective masses (left) and areas (right) as a function of true neutrino energy for a generic toy detector with fiducial mass of 10 Mt. The dependency of the effective masses on energy is given in Equation (A.1). Cross sections are from Equation (A.3), except for  $\nu_\tau$  and  $\bar{\nu}_\tau$  interactions, which are interpolated from [46]. Effective masses are the same for neutrinos and anti-neutrinos. See text for details.

$x \rightarrow (x - \mu')/\sigma'$ . These parameters again only depend on  $E_{\text{true}}$ . The CC distributions are assumed identical for all flavors, and are shown in Figure A.16:

$$\mu'_{\Delta E_\nu}{}^{\text{CC}}(E_{\text{true}}) = 0, \quad \sigma'_{\Delta E_\nu}{}^{\text{CC}}(E_{\text{true}}) = \left( \frac{0.4}{\sqrt{E_{\text{true}}/\text{GeV}}} + 0.1 \right) \times E_{\text{true}}. \quad (\text{A.7})$$

For NC interactions, we take a spread that scales with  $E_{\text{true}}$  in the same way  $\sigma'_{\Delta E_\nu}{}^{\text{CC}}$  does, but assume a non-zero mean due to the undetected energy carried away by the outgoing neutrino:  $\mu'_{\Delta E_\nu}{}^{\text{NC}}(E_{\text{true}}) = -0.6E_{\text{true}}$ .

Note that each energy and cosine zenith residual distribution is renormalized such that its integral over the physical region ( $\Delta E_\nu + E_{\text{true}} \geq 0$  and  $-1 \leq (\Delta \cos \vartheta + \cos \vartheta_{\text{true}}) \leq 1$ ) yields 1.

### Appendix A.3. Event Classification

Correctly identifying few-GeV CC muon neutrino interactions with relatively sparsely instrumented neutrino telescopes in water/ice is difficult mainly for two reasons. The track length of a near minimum ionizing muon is only on the order of a few meters, comparable to the extent of an electromagnetic cascade of the same energy. Also, photon scattering lengths similar to the horizontal spacing between photomultiplier tubes smear out the Cherenkov ring around the muon direction, which is otherwise observed at a specific angle with respect to the muon direction in the medium.

We take into account the muon neutrino CC (“track”) identification efficiency  $p_{\text{track}}^{\mu, \text{CC}}$  improving with (reconstructed) neutrino energy,  $E_{\text{reco}}$ , by setting

$$p_{\text{track}}^{\mu, \text{CC}} \equiv p_{\text{track}}^{\mu, \text{CC}}(E_{\text{reco}}) = 0.9 \times \left( 1 - e^{-0.2 \times (E_{\text{reco}}/\text{GeV} + 0.6)} \right). \quad (\text{A.8})$$

This reflects the track length of the secondary muon increasing linearly with its energy, but also the possible production of a low-energy muon which cannot be distinguished from the

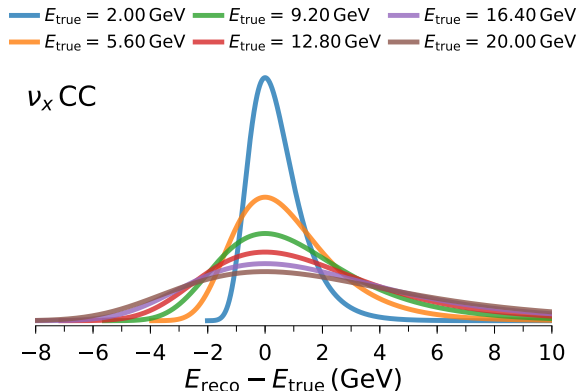


Figure A.16: Example energy resolution functions (right-skewed Gumbel) used for all CC interactions, as given by Equation (A.7). The modes of the corresponding NC resolution functions are shifted by  $-0.6E_{\text{true}}$  with respect to the distributions depicted here.

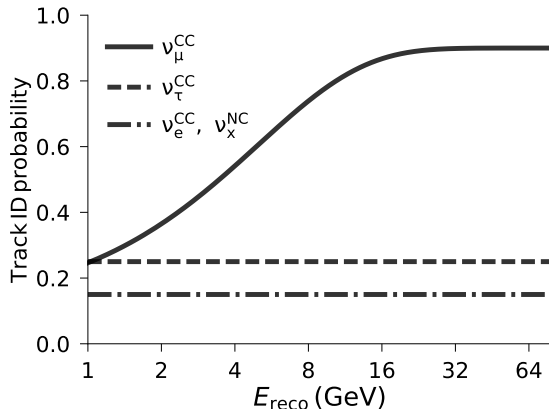


Figure A.17: Event classification efficiencies implemented as functions of reconstructed neutrino energy. Shown is the probability to identify an event of a given type as “track-like”. Events are identified as “cascade-like” with complementary probabilities.

accompanying hadronic cascade even for higher-energy muon neutrino CC interactions. All other (in)efficiencies are assumed to be constant:

$$p_{\text{track}}^{e,\text{CC}}(E_{\text{reco}}) = p_{\text{track}}^{\text{NC}}(E_{\text{reco}}) = 0.15, \quad (\text{A.9})$$

$$p_{\text{track}}^{\tau,\text{CC}}(E_{\text{reco}}) = 0.25. \quad (\text{A.10})$$

These are shown in Figure A.17. The probability to identify any event as “cascade-like” for a given reconstructed energy is just the complementary probability to that of the track identification.

When a toy MC event is subject to this classification, we assign it one of two discrete numbers—representative of either identification as track or cascade—with the above probabilities.

## Appendix B. Uncertainty in Significance

Under the assumption that the test statistic  $\mathcal{T}$  under two hypotheses  $H_1$  and  $H_2$  is normally distributed (with means  $\mu_1$  and  $\mu_2$  and with identical standard deviation  $\sigma$ ), the number of standard deviations ( $n_\sigma$ ) separating the two hypotheses can be written as  $n_\sigma = |\mu_1 - \mu_2|/\sigma$  (corresponding to a one-sided hypothesis test and a one-sided conversion from p-value). Sampling each distribution with  $N_p$  pseudo-experiments results in the following uncertainties for mean and standard deviation (see for example [47])

$$\Delta\mu = \frac{\sigma}{\sqrt{N_p}}, \quad (\text{B.1})$$

$$\Delta\sigma = \frac{\sigma}{\sqrt{2(N_p - 1)}}. \quad (\text{B.2})$$

Since the combination of the quantities is linear, we can perform simple error propagation, so that the relative uncertainty in significance becomes (with  $\oplus$  denoting sum in quadrature)

$$\frac{\Delta n_\sigma}{n_\sigma} = \frac{\Delta \sigma}{\sigma} \oplus \frac{\Delta |\mu_1 - \mu_2|}{|\mu_1 - \mu_2|}. \quad (\text{B.3})$$

Using

$$\Delta |\mu_1 - \mu_2| = \Delta \mu_1 \oplus \Delta \mu_2 = \sqrt{\frac{2}{N_p}} \sigma \quad (\text{B.4})$$

the second term simplifies to

$$\frac{\Delta |\mu_1 - \mu_2|}{|\mu_1 - \mu_2|} = \sqrt{\frac{2}{N_p}} \frac{\sigma}{|\mu_1 - \mu_2|} = \sqrt{\frac{2}{N_p}} \frac{1}{n_\sigma}. \quad (\text{B.5})$$

Substituting Equations (B.5) and (B.1) into Equation (B.3) yields

$$\frac{\Delta n_\sigma}{n_\sigma} = \frac{1}{\sqrt{2(N_p - 1)}} \oplus \sqrt{\frac{2}{N_p n_\sigma^2}} = \sqrt{\frac{1}{2(N_p - 1)} + \frac{2}{N_p n_\sigma^2}}. \quad (\text{B.6})$$

The absolute error on the number of standard deviations and its approximation for large  $N_p$  then follow immediately as

$$\Delta n_\sigma = \sqrt{\frac{n_\sigma^2}{2(N_p - 1)} + \frac{2}{N_p}} \xrightarrow{(N_p \gg 1)} \frac{1}{\sqrt{N_p}} \sqrt{\frac{n_\sigma^2}{2} + 2}. \quad (\text{B.7})$$

## References

- [1] A. Ishihara, The IceCube Upgrade - Design and Science Goals, in: 36th International Cosmic Ray Conference (ICRC 2019) Madison, Wisconsin, USA, July 24-August 1, 2019, 2019. [arXiv:1907.11699](#).
- [2] M. G. Aartsen, et al., PINGU: a vision for neutrino and particle physics at the South Pole, *J. Phys. G44* (5) (2017) 054006. [arXiv:1607.02671](#), [doi:10.1088/1361-6471/44/5/054006](#).
- [3] M. G. Aartsen, et al., Letter of intent: the Precision IceCube Next Generation Upgrade (PINGU), (2017). [arXiv:1401.2046v2](#).
- [4] S. Adrian-Martinez, et al., Letter of intent for KM3NeT 2.0, *J. Phys. G43* (8) (2016) 084001. [arXiv:1601.07459](#), [doi:10.1088/0954-3899/43/8/084001](#).
- [5] T. Glüsenkamp, A unified perspective on modified Poisson likelihoods for limited Monte Carlo data. [arXiv:1902.08831](#).
- [6] C. A. Argüelles, A. Schneider, T. Yuan, A binned likelihood for stochastic models, *JHEP* 06 (2019) 030. [arXiv:1901.04645](#), [doi:10.1007/JHEP06\(2019\)030](#).
- [7] R. Barlow, Introduction to statistical issues in particle physics, Statistical problems in particle physics, astrophysics and cosmology. Proceedings, Conference, PHYSTAT 2003, Stanford, USA, September 8-11, 2003, C030908 (2003) MOAT002. [arXiv:physics/0311105](#).
- [8] G. Cowan, *Statistical Data Analysis*, Oxford University Press, 1998.
- [9] R. H. Byrd, P. Lu, J. Nocedal, A limited memory algorithm for bound constrained optimization, *SIAM J. Sci. Comput.* 16 (1995) 1190–1208.

- [10] G. Cowan, K. Cranmer, E. Gross, O. Vitells, Asymptotic formulae for likelihood-based tests of new physics, *Eur. Phys. J. C* 71 (2011) 1554, [Erratum: *Eur. Phys. J. C* 73,2501(2013)]. [arXiv:1007.1727](#), [doi:10.1140/epjc/s10052-011-1554-0](#), [10.1140/epjc/s10052-013-2501-z](#).
- [11] M. Blennow, P. Coloma, P. Huber, T. Schwetz, Quantifying the sensitivity of oscillation experiments to the neutrino mass ordering, *JHEP* 2014 (3) (2014) 28. [doi:10.1007/JHEP03\(2014\)028](#).
- [12] K. S. Cranmer, Kernel estimation in high-energy physics, *Comput. Phys. Commun.* 136 (2001) 198–207. [arXiv:hep-ex/0011057](#), [doi:10.1016/S0010-4655\(00\)00243-5](#).
- [13] D. W. Scott, On optimal and data-based histograms, *Biometrika* 66 (3) (1979) 605. [doi:10.1093/biomet/66.3.605](#).
- [14] Y. Fukuda, et al., Evidence for oscillation of atmospheric neutrinos, *Phys. Rev. Lett.* 81 (1998) 1562–1567. [arXiv:hep-ex/9807003](#), [doi:10.1103/PhysRevLett.81.1562](#).
- [15] Q. R. Ahmad, et al., Measurement of the rate of  $\nu_e + d \rightarrow p + p + e^-$  interactions produced by  $^8B$  solar neutrinos at the Sudbury Neutrino Observatory, *Phys. Rev. Lett.* 87 (2001) 071301. [arXiv:nucl-ex/0106015](#), [doi:10.1103/PhysRevLett.87.071301](#).
- [16] B. Aharmim, et al., Combined analysis of all three phases of solar neutrino data from the Sudbury Neutrino Observatory, *Phys. Rev. C* 88 (2013) 025501. [arXiv:1109.0763](#), [doi:10.1103/PhysRevC.88.025501](#).
- [17] F. Capozzi, E. Lisi, A. Marrone, A. Palazzo, Current unknowns in the three neutrino framework, *Prog. Part. Nucl. Phys.* 102 (2018) 48–72. [arXiv:1804.09678](#), [doi:10.1016/j.pnpnp.2018.05.005](#).
- [18] P. F. De Salas, S. Gariazzo, O. Mena, C. A. Ternes, M. Tórtola, Neutrino Mass Ordering from Oscillations and Beyond: 2018 Status and Future Prospects, *Front. Astron. Space Sci.* 5 (2018) 36. [arXiv:1806.11051](#), [doi:10.3389/fspas.2018.00036](#).
- [19] NuFIT 4.0 (2018), [www.nu-fit.org](#).
- [20] I. Esteban, M. C. Gonzalez-Garcia, A. Hernandez-Cabezudo, M. Maltoni, T. Schwetz, Global analysis of three-flavour neutrino oscillations: synergies and tensions in the determination of  $\theta_{23}$ ,  $\delta_{CP}$ , and the mass ordering, *JHEP* 01 (2019) 106. [arXiv:1811.05487](#), [doi:10.1007/JHEP01\(2019\)106](#).
- [21] C. Patrignani, et al., Review of particle physics, *Chin. Phys.* C40 (10) (2016) 100001. [doi:10.1088/1674-1137/40/10/100001](#).
- [22] L. Wolfenstein, Neutrino oscillations in matter, *Phys. Rev. D* 17 (1978) 2369.
- [23] S. P. Mikheev, A. Y. Smirnov, Resonant amplification of neutrino oscillations in matter and solar neutrino spectroscopy, *Nuovo Cim.* C9 (1986) 17–26.
- [24] S. T. Petcov, S. Toshev, Three neutrino oscillations in matter: analytical results in the adiabatic approximation, *Phys. Lett. B* 187 (1987) 120–126. [doi:10.1016/0370-2693\(87\)90083-9](#).
- [25] E. K. Akhmedov, M. Maltoni, A. Yu. Smirnov, Oscillations of high energy neutrinos in matter: precise formalism and parametric resonance, *Phys. Rev. Lett.* 95 (2005) 211801. [arXiv:hep-ph/0506064](#), [doi:10.1103/PhysRevLett.95.211801](#).
- [26] K. Abe, et al., Letter of intent: the Hyper-Kamiokande experiment — detector design and physics potential —. [arXiv:1109.3262](#).
- [27] J. Neyman, E. S. Pearson, On the problem of the most efficient tests of statistical hypotheses, *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences* 231 (694-706) (1933) 289–337. [doi:10.1098/rsta.1933.0009](#).
- [28] M. C. Gonzalez-Garcia, M. Maltoni, T. Schwetz, Updated fit to three neutrino mixing: status of leptonic CP violation, *JHEP* 2014 (11) (2014) 52. [doi:10.1007/JHEP11\(2014\)052](#).
- [29] NuFIT 2.0 (2014), [www.nu-fit.org](#).
- [30] M. Honda, M. Sajjad Athar, T. Kajita, K. Kasahara, S. Midorikawa, Atmospheric neutrino flux calculation using the NRLMSISE-00 atmospheric model, *Phys. Rev. D* 92 (2) (2015) 023004. [arXiv:1502.03916](#), [doi:10.1103/PhysRevD.92.023004](#).
- [31] G. D. Barr, T. K. Gaisser, S. Robbins, T. Stanev, Uncertainties in atmospheric neutrino fluxes, *Phys. Rev. D* 74 (2006) 094009. [arXiv:astro-ph/0611266](#), [doi:10.1103/PhysRevD.74.094009](#).
- [32] J. Evans, D. G. Gamez, S. D. Porzio, S. Söldner-Rembold, S. Wren, Uncertainties in atmospheric muon-neutrino fluxes arising from cosmic-ray primaries, *Phys. Rev. D* 95 (2) (2017) 023012. [arXiv:1612.03219](#),

- doi:10.1103/PhysRevD.95.023012.
- [33] V. Barger, K. Whisnant, S. Pakvasa, R. J. N. Phillips, Matter effects on three-neutrino oscillations, *Phys. Rev. D* 22 (1980) 2718–2726. doi:10.1103/PhysRevD.22.2718.
  - [34] R. Wendell, Prob3++ software for computing three flavor neutrino oscillation probabilities, <http://www.phy.duke.edu/~raw22/public/Prob3++/>, 2012.
  - [35] A. M. Dziewonski, D. L. Anderson, Preliminary reference Earth model, *Physics of the Earth and planetary interiors* 25 (4) (1981) 297 – 356. doi:http://dx.doi.org/10.1016/0031-9201(81)90046-7.
  - [36] A. Bulmahn, M. H. Reno, Secondary atmospheric tau neutrino production, *Phys. Rev. D* 82 (2010) 057302. arXiv:1007.4989, doi:10.1103/PhysRevD.82.057302.
  - [37] S. Wren, Neutrino Mass Ordering Studies with IceCube-DeepCore, Ph.D. thesis, Manchester U. (2018). URL [https://www.research.manchester.ac.uk/portal/en/theses/neutrino-mass-ordering-studies-with-icecubedeepcore\(70414fde-3bef-4028-877b-5e1e86b2165d\).html](https://www.research.manchester.ac.uk/portal/en/theses/neutrino-mass-ordering-studies-with-icecubedeepcore(70414fde-3bef-4028-877b-5e1e86b2165d).html)
  - [38] R. G. Calland, A. C. Kaboth, D. Payne, Accelerated event-by-event neutrino oscillation reweighting with matter effects on a GPU, *JINST* 9 (2014) P04016. arXiv:1311.7579, doi:10.1088/1748-0221/9/04/P04016.
  - [39] J. Nickolls, I. Buck, M. Garland, K. Skadron, Scalable parallel programming with cuda, 2008 IEEE Hot Chips 20 Symposium (HCS) (2008) 1–2.
  - [40] M. Wallraff, C. Wiebusch, Calculation of oscillation probabilities of atmospheric neutrinos using nuCraft, *Comput. Phys. Commun.* 197 (2015) 185–189. arXiv:1409.1387, doi:10.1016/j.cpc.2015.07.010.
  - [41] T. K. Gaisser, T. Stanev, Path length distributions of atmospheric neutrinos, *Phys. Rev. D* 57 (1998) 1977–1982. doi:10.1103/PhysRevD.57.1977.
  - [42] I. S. Abramson, On bandwidth variation in kernel estimates—a square root law, *Ann. Statist.* 10 (4) (1982) 1217–1223. doi:10.1214/aos/1176345986.
  - [43] Z. I. Botev, J. F. Grotowski, D. P. Kroese, Kernel density estimation via diffusion, *Ann. Statist.* 38 (5) (2010) 2916–2957. doi:10.1214/10-AOS799.
  - [44] S. J. Sheather, M. C. Jones, A reliable data-based bandwidth selection method for kernel density estimation, *Journal of the Royal Statistical Society. Series B (Methodological)* 53 (3) (1991) 683–690.
  - [45] J. A. Formaggio, G. P. Zeller, From eV to EeV: neutrino cross sections across energy scales, *Rev. Mod. Phys.* 84 (2012) 1307–1341. arXiv:1305.7513, doi:10.1103/RevModPhys.84.1307.
  - [46] A. Gazizov, M. Kowalski, K. S. Kuzmin, V. A. Naumov, C. Spiering, Neutrino-nucleon cross sections at energies of megaton-scale detectors, *EPJ Web Conf.* 116 (2016) 08003. arXiv:1604.02092, doi:10.1051/epjconf/201611608003.
  - [47] D. Sivia, *Data Analysis: A Bayesian Tutorial*, Oxford University Press, 2006.