

Development of Bayesian analysis program for extraction of polarisation observables at CLAS

S. Lewis, D. Ireland, W. Vanderbauwhede

SUPA, School of Physics and Astronomy, University of Glasgow, Glasgow, G12 8QQ, United Kingdom

E-mail: `s.lewis.2@research.gla.ac.uk`

Abstract. At the mass scale of a proton, the strong force is not well understood. Various quark models exist, but it is important to determine which quark model(s) are most accurate. Experimentally, finding resonances predicted by some models and not others would give valuable insight into this fundamental interaction. Several labs around the world use photoproduction experiments to find these missing resonances. The aim of this work is to develop a robust Bayesian data analysis program for extracting polarisation observables from pseudoscalar meson photoproduction experiments using CLAS at Jefferson Lab. This method, known as nested sampling, has been compared to traditional methods and has incorporated data parallelisation and GPU programming. It involves an event-by-event likelihood function, which has no associated loss of information from histogram binning, and results can be easily constrained to the physical region. One of the most important advantages of the nested sampling approach is that data from different experiments can be combined and analysed simultaneously. Results on both simulated and previously analysed experimental data for the $K^+\Lambda$ channel will be discussed.

1. Introduction

The strong force at the level of the mass of the nucleon is not well understood. Experiments addressing this fundamental problem of nuclear physics have been performed by several international collaborations, including CLAS [1] at Jefferson Lab and A2 at Mainz. Pseudoscalar meson photoproduction reactions can be used to learn about the excited nucleon spectrum. The use of Bayesian analysis in hadronic physics has several benefits over more conventional statistical methods. Nested sampling [2] is an algorithm based on Bayesian statistics that can extract more information from data than some other conventional methods. Data parallelism and GPU programming can be implemented to speed up the algorithm, and there are several tools that can be used to do this.

2. Background Physics Theory

Quantum chromodynamics explains the strong interaction in general, but cannot be solved at the mass of the proton. Many quark models exist that attempt to address this issue [3]. Quark models can be supported or excluded based on experimentally finding mass resonances that the models predict. One way to find these resonances is through pseudoscalar meson photoproduction [1] reactions, where a photon beam is incident on a stationary nucleon target.

Pseudoscalar meson photoproduction reactions are completely described by four complex amplitudes, a_1 , a_2 , a_3 and a_4 [4]. These amplitudes can be accessed experimentally through 15 polarisation observables, listed in Table 1, in addition to the differential cross-section. Polarising different aspects of the experiment provides access to a different subset of observables, and a well-chosen set of observables are required to make a complete measurement. There are three single-spin observables (Σ , P , T), four double-spin beam-recoil observables (C_x , C_z , O_x , O_z), four double-spin beam-target observables (E , F , G , H), and four double-spin target-recoil observables (T_x , T_z , L_x , L_z). These observables are bilinear combinations of the amplitudes, and measuring a subset of observables should constrain the others. In other words, information can be gained about observables that are not directly measured through their correlations.

Table 1. Observables in terms of Complex Amplitudes [4]

Observable	Type	Amplitude Combination
Σ	Single	$ a_1 ^2 + a_2 ^2 - a_3 ^2 - a_4 ^2$
P		$ a_1 ^2 - a_2 ^2 + a_3 ^2 - a_4 ^2$
T		$ a_1 ^2 - a_2 ^2 - a_3 ^2 + a_4 ^2$
E	Beam-target	$2\Re(a_1 a_3^* + a_2 a_4^*)$
F		$2\Im(a_1 a_3^* - a_2 a_4^*)$
G		$2\Im(a_1 a_3^* + a_2 a_4^*)$
H		$-2\Re(a_1 a_3^* - a_2 a_4^*)$
C_x	Beam-recoil	$-2\Im(a_1 a_4^* - a_2 a_3^*)$
C_z		$2\Re(a_1 a_4^* + a_2 a_3^*)$
O_x		$2\Re(a_1 a_4^* - a_2 a_3^*)$
O_z		$2\Im(a_1 a_4^* + a_2 a_3^*)$
T_x	Target-recoil	$2\Re(a_1 a_2^* - a_3 a_4^*)$
T_z		$2\Im(a_1 a_2^* - a_3 a_4^*)$
L_x		$-2\Im(a_1 a_2^* + a_3 a_4^*)$
L_z		$2\Re(a_1 a_2^* + a_3 a_4^*)$

In this work, data from an experiment involving a linearly polarised photon beam and an unpolarised proton (liquid hydrogen) target were analysed. The reaction is shown in Equation 1.

$$\gamma p \rightarrow K^+ \Lambda \rightarrow K^+ p \pi^- \quad (1)$$

The resulting Λ then decays to a proton and pion. The scattered kaon and recoiling proton were detected, and the polarisation of the Λ was determined through the self-analysing properties of the hyperon. This provided access to the observables Σ , P , T , O_x and O_z .

The experimental data was collected by the CLAS Collaboration, based in Hall B at the Thomas Jefferson National Accelerator Facility. Hall B houses the CLAS (CEBAF Large Acceptance Spectrometer) detector [5], which has a wide solid angle acceptance and is well understood.

3. Analysis Method

The current most popular method of analysis involves the classical χ^2 binned fitting approach. This method involves binning data into histograms and calculating polarisation observables based on asymmetries of these binned histograms [7]. There are several disadvantages to this approach: Information is lost through histogram binning, observables are not bound to the physical region, and observables are treated independently and as free variables (the information

yield is not maximised; correlations between observables are not used). In order to address some of these problems, a new analysis method was explored using Bayesian statistics.

Bayesian analysis has recently become a popular alternative to conventional analysis methods. A distinguishing feature of the Bayesian approach is its use of a prior distribution, which can be thought of as an initial approximation to the results, before any data is taken into consideration. The prior can contain any physical constraints known about the system, and a good choice of prior can often speed up the analysis [8]. A likelihood function is then applied to the prior distribution, and a posterior distribution is produced (the statistics of which give the results). In this work, the nested sampling algorithm was applied to the analysis of data from pseudoscalar meson photoproduction experiments, and provides certain interesting features over the more conventional analysis methods currently being used. The fundamental difference between this Bayesian approach and the χ^2 binned fitting method is that the amplitude space is explored, and the observables are not completely independent variables. Observables are calculated from the complex amplitudes, using the equations listed in Table 1. The use of an event-by-event likelihood function rather than binned histograms reduces the amount of information lost, and the prior distribution includes conditions that constrain the observables to the physical region. These features allow the information extracted from the data to be maximised. One problem was found with the nested sampling program, however, as the run-time was significantly longer than the χ^2 program. This led to the investigation into data parallelism and GPU programming techniques, discussed in Section 4.

The nested sampling analysis program was benchmarked against the χ^2 fitting method, using 10 sets of simulated data, and three dataset sizes. The difference between the input observable values and the extracted values were compared for each of the two methods, and the resulting plot is shown in Figure 1 (top). This plot indicates that the nested sampling approach produced more accurate values, even when the statistics were poor. The second plot in Figure 1 shows an measure of the understanding of the errors, where the optimal value would be expected to equal the number of observables being extracted. Here, the ideal value would have the square root of the sum of the residuals to be 5. It is seen that nested sampling achieves this value more consistently than does the binned χ^2 fit.

A further benefit of the nested sampling approach is the potential to combine data from different experiments and extract observables that inherently obey the constraints imposed by their correlations. This has been tested using simulated data and shown to be successful, but has yet to be attempted using data from real experiments.

4. Data Parallelism

The exceedingly long run-time of the nested sampling program spurred efforts into optimisation and parallelism. The code, written in C++ (and using the ROOT framework), was optimised before any parallelism was explored. It was determined that the event-by-event likelihood function was the primary bottleneck. Implementations were developed using OpenCL (Open Computing Language) [9] and OpenMP (Open Multiprocessing)[10]. OpenCL is an open-source framework that allows code to be parallelised and run on the GPU or the CPU, but is dependent on the hardware. Code written for an AMD GPU, for example, would need to be fine-tuned to obtain the best results on an NVIDIA GPU. To parallelise a program using OpenCL, the code must consist of two parts: the host code and the kernel code. The kernel code contains the function or algorithm that is parallelised, and is called by the host program. There is a substantial amount of “boiler plate” code that is required to initialise the OpenCL environment, which can make the programming complicated and intimidating. Collaboration with the School of Computing Science led to the use of OclWrapper (<https://github.com/wimvanderbauwhede/OpenCLIntegration>), a wrapper class developed by W. Vanderbauwhede that greatly simplifies this task. An OpenCL implementation of the nested

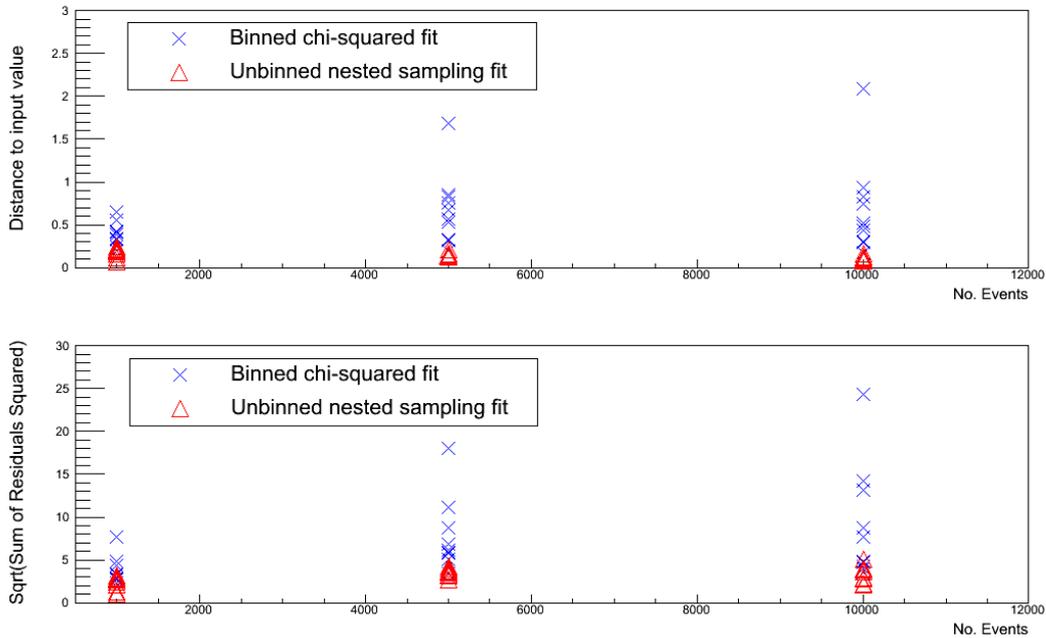


Figure 1. Benchmarking results comparing the accuracy of the nested sampling method with the binned χ^2 approach.

sampling program [11] was developed and tested on a GPU as well as a multicore CPU. The other option explored was an implementation using OpenMP, an open-source multithreading API. OpenMP uses a series of compiler directives to allow programmers to manipulate threads and parallelise algorithms. This was a much simpler option, requiring only a few additional lines of code inserted into the existing program, and is supported by most recent compilers. The OpenMP implementation of the program was limited to running on the CPU only. The various implementations (including the optimised, unthreaded and unparallelised version) were timed for three different dataset sizes, and the results are shown in Figure 2.

It was found that for the smaller datasets, OpenMP provided the fastest run-time, and in fact, OpenCL run on the GPU provided the slowest (slower even than the unparallelised code). This is due to the difference in data transfer times related to the hardware devices. The time required to transfer data from the main memory to the GPU memory far exceeds that required to transfer data to the CPU. For the OpenCL implementations, the overhead caused by transferring data exceeded any speed-up in the calculations. As the dataset size increases, it can be seen that this speed-up from calculations begins to catch up with the data transfer overhead, and it is expected that if greater datasets were tested, the GPU implementation would provide the best speed-up. For the specific physics problem faced here, however, the datasets that are anticipated are more towards the lower number of events, often a few 10^3 .

5. Results

Previously analysed data from a CLAS experiment was used to test the nested sampling program against conventional methods. Whilst nested sampling does not bin data into histograms in order to perform its analysis, the datasets that are analysed are binned by photon energy (W) and scattering angle in the centre of mass frame ($\cos(\theta_K^{CM})$).

The resulting observables extracted via the nested sampling analysis program were compared

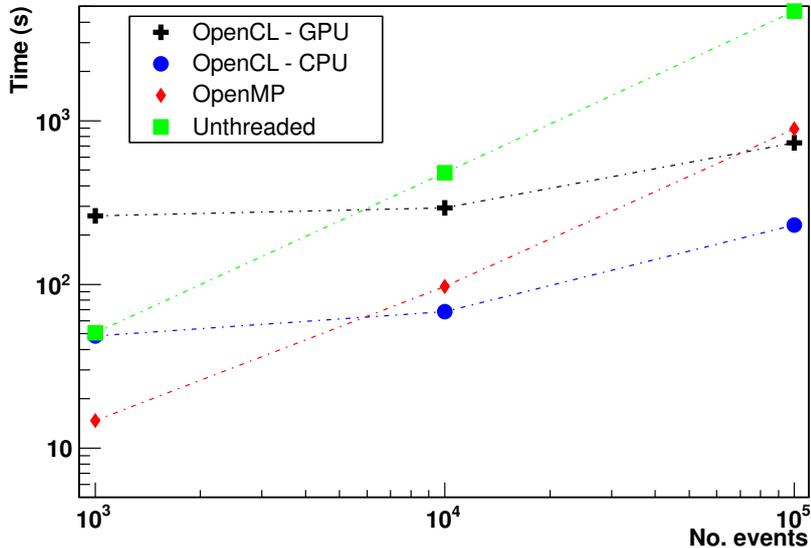


Figure 2. Run-times of each implementation with increasing dataset size.

to results found in previous analyses [6] of the same data. Figure 3 shows the comparison plot for one observable, Σ . It can be seen that the nested sampling method produces values that are sufficiently close and consistent with previous results.

One of the benefits of using nested sampling becomes apparent with one type of output. By exploring the amplitude space, information about observables not being directly measured is discovered and can be shown as two-dimensional likelihood plots. Figure 4 shows a series of plots of all 15 observables in one energy bin. Each blob, or band, shows the resulting posterior distribution for each of the 15 observables, each corresponding to one angular bin. The plot shows how much the other observables are constrained by the information found from the data. This is information that was previously not obtainable and makes a significant difference in the impact of these pseudoscalar meson photoproduction experiments.

6. Summary

In order to address the problem of understanding the strong interaction on the hadronic level, it is important to find missing resonances that are predicted by quark models, as yet unobserved. This can be done through pseudoscalar meson photoproduction experiments, such as those performed by the CLAS Collaboration at Jefferson Lab. Pseudoscalar meson photoproduction reactions can be described completely by four complex amplitudes, which can be accessed experimentally through polarisation observables. These observables are bilinear combinations of the amplitudes, suggesting that measuring some observables will inherently provide information about all the others. Current methods of analysis treat observables as independent, free variables and therefore do not maximise the information extracted from the data. A new Bayesian approach called nested sampling has been developed, whose distinguishing feature is the use of amplitude space and preserve the correlations between polarisation observables. Unfortunately, this method has a significantly longer run-time than the more conventional analysis methods. Data parallelism and multithreading can be used to speed up the program quite effectively, and studies were done into the most efficient solution. It was found that introducing multithreading using OpenMP gave the fastest run-time at the dataset sizes that will likely be faced in analysing

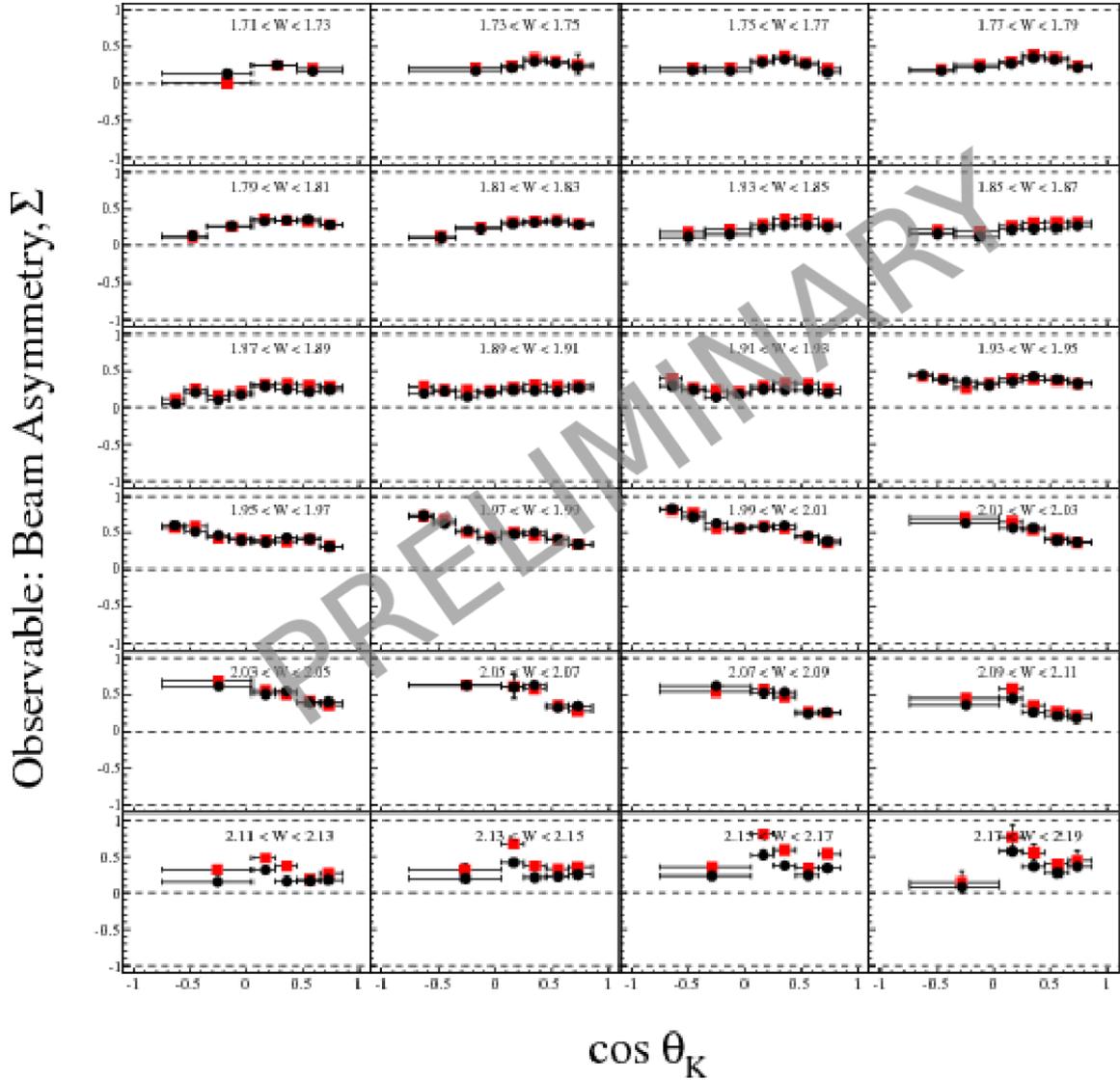
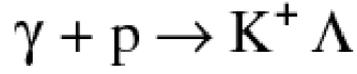


Figure 3. Preliminary results from CLAS data (colour online): Previously obtained results (red squares) [7] are compared to the results from the nested sampling program (black circles).

pseudoscalar meson photoproduction experiments. Future work will be focused on analysing combined datasets of multiple experiments and extracting observables in a way that retains the constraints of their correlations.

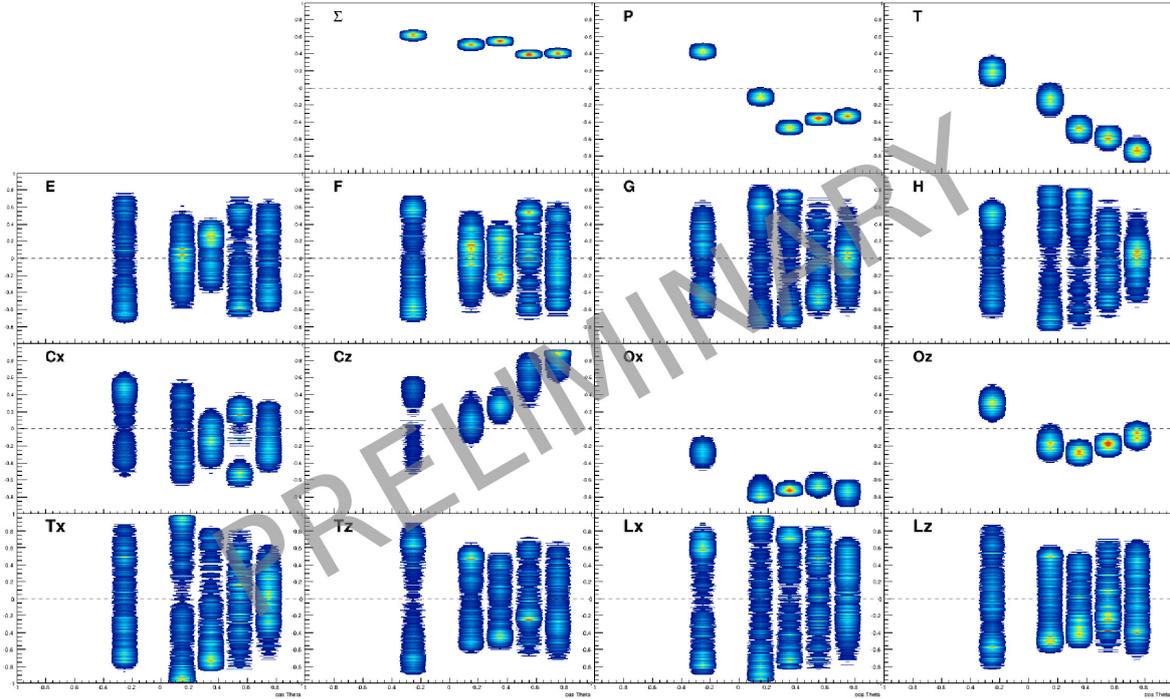


Figure 4. Plot showing the posterior distributions of all 15 observables in one energy bin, $2.03 < W < 2.05$ GeV (colour online). For each cosine bin, the resulting distribution is shown, where brighter areas suggest the more likely values of the observable.

References

- [1] I. Aznauryan, V. Burkert, T.-S. Lee, V. Mokeev, arXiv:1102.0597v3 [nucl-ex], February 2011
- [2] J. Skilling, *Bayesian Analysis* 1 4, 833 (2006)
- [3] E. Klempt, J.-M. Richard, *Rev. Mod. Phys.* 82, 1095-1153 (2010)
- [4] D. G. Ireland, *Phys. Rev. C* 82, 025204 (2010)
- [5] B. A. Mecking et al, *Nucl. Instrum. Methos Phys. Res., Sect. A* 503, 513 (2003)
- [6] D. G. Ireland, Private Communication.
- [7] C. A. Paterson, *Polarization Observables in Strangeness Photoproduction with CLAS at Jefferson Lab*. PhD thesis, University of Glasgow (2008)
- [8] D. Sivia and J. Skilling, *Data Analysis - A Bayesian Tutorial*. 2nd ed. Oxford Science Publications (2006)
- [9] Khronos Group, 2011. *The OpenCL Specification, Version 1.2* [online] Available at: <http://www.khronos.org> [Accessed November 2011]
- [10] OpenMP (2011) *OpenMP Application Program Interface* [online] Available at: <http://openmp.org/wp/openmp-specifications/> [Accessed February 2012]
- [11] W. Vanderbauwhede, S. Lewis, D. G. Ireland, *Implementing Data Parallelisation in a Nested-Sampling Monte Carlo Algorithm*, tbp in Proc. HPCS'13, The 2013 International Conference on High Performance Computing & Simulation, Helsinki, Finland, July 2013