

Rheinisch-Westfälische Technische Hochschule Aachen

## **Master Thesis**

# **Denoising Methods for Multi-Dimensional Photoemission Spectroscopy**

SUBMITTED BY

**Muhammad Zain Sohail**

Rheinisch-Westfälische Technische Hochschule Aachen  
Christian-Albrechts-Universität zu Kiel  
Deutsches Elektronen Synchrotron

SUPERVISORS

**Prof. Dr. Benjamin Berkels**

Rheinisch-Westfälische Technische Hochschule Aachen

**Prof. Dr. Kai Rossnagel**

Christian-Albrechts-Universität zu Kiel  
Deutsches Elektronen Synchrotron



Aachen, October 2024

# Summary

In the realm of photoemission spectroscopy, the exploration of large multi-dimensional phase spaces necessitates time-intensive data acquisition to ensure statistical robustness. Despite the unparalleled capabilities of free-electron lasers (FELs), in peak brightness and ultra- short pulsed X-rays, the limitations of low repetition rates prolong the data acquisition process. This impedes the agility of decision making that could otherwise enhance experimental results in the limited and valuable beam-time. By employing denoising strategies to mitigate noise while preserving intrinsic information, our proposed approach aims to streamline the data acquisition process, and effectively manage the escalating size and complexity of multi-dimensional photoemission data.

# Contents

<b>Summary</b>	<b>ii</b>
<b>Contents</b>	<b>iii</b>
<b>List of Acronyms</b>	<b>v</b>
<b>Glossary</b>	<b>vi</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Photoemission Spectroscopy</b>	<b>2</b>
2.1 Light-Matter Interaction . . . . .	2
2.2 Light Sources . . . . .	2
2.3 SASE FELs . . . . .	3
2.4 HEXTOF Setup at FLASH . . . . .	3
<b>3 Imaging</b>	<b>5</b>
3.1 Binning as means of Imaging . . . . .	5
3.2 SASE fluctuations . . . . .	5
3.3 Describing the data GdW, WSe2, GrIr, and new one . . . . .	5
3.4 Creating the dataset Corrections Calibrations etc . . . . .	5
<b>4 Characterizing Photon and Photoelectron Statistics</b>	<b>6</b>
4.1 Poisson Process assumption . . . . .	6
4.1.1 Before filtering . . . . .	7
4.1.2 After filtering . . . . .	7
4.2 Simulate Noise . . . . .	7
4.3 Statistical testing . . . . .	7
4.4 Chi-squared Goodness of Fit Test . . . . .	11
<b>5 Image Reconstruction</b>	<b>12</b>
5.1 Address reconstruction/denoising schemes . . . . .	12

5.2	BM3D: Denoising in sparse domain . . . . .	12
5.2.1	Anscombe: Variance Stabilization Transform . . . . .	12
5.3	Statistical Learning . . . . .	12
5.3.1	Optimal Loss Function . . . . .	13
5.3.2	Regularization . . . . .	13
5.3.3	Optimization . . . . .	14
5.4	Noise2Noise: Deep Learning framework . . . . .	14
5.4.1	Convolutional Neural Networks . . . . .	14
5.4.2	Autoencoder . . . . .	14
<b>6</b>	<b>Conclusion and Outlook</b>	<b>15</b>
	<b>Bibliography</b>	<b>16</b>
	<b>Acknowledgements</b>	<b>17</b>
	<b>Contributions</b>	<b>18</b>
	<b>Declaration</b>	<b>19</b>
	<b>Appendices</b>	<b>20</b>
<b>A</b>	<b>Transforming Raw Data to structured format</b>	<b>21</b>
<b>B</b>	<b>Mathematical Background</b>	<b>24</b>
B.1	Measure Space and Measures . . . . .	24
B.2	Probability . . . . .	25
B.3	Landau Notation . . . . .	25
B.4	Statistical Inference . . . . .	26
B.4.1	Testing/Confidence Intervals . . . . .	26
B.4.2	Optimal parameter estimation . . . . .	26
<b>C</b>	<b>Deep Learning</b>	<b>27</b>
C.1	Infrastructure . . . . .	27

# List of Acronyms

**BAM** beam arrival monitor. 19

**CDF** cumulative distribution function. 11

**DESY** Deutsches Elektronen-Synchrotron. 3

**DLD** delay line detector. 19

**FEL** free-electron laser. v

**GMD** gas monitor detector. 19

**HDF5** hierarchical data format version 5. 19

**OpenCOMPES** open community of multidimensional photoemission spectroscopy.  
19

**PES** photoemission/photoelectron spectroscopy. 2

**SASE** self-amplified spontaneous emission. v

**SED** Single Event DataFrame. 19

# Glossary

**Beamline** A path leading the photons from the particle accelerator to the experimental end-station.. 19

**Free-Electron Laser** is an x-ray radiation source; fundamentally comprising of a linear particle accelerator and an Undulator (or a series of undulators). The accelerator produces a bunched electron beam similar to that of a synchrotron, which can be compressed to reach ultrashort pulse duration (femtosecond) with peak brightness many orders of magnitude above synchrotrons. Since the electrons move in a vacuum, it is termed Free-Electron in comparison to traditional lasers which are bound by the materials energy levels. Whereas, it is called a Laser due to there being light amplification and the shared properties with traditional optical lasers such as high pulse energy and being coherent. Taken from [7]. 1

**Microbunch** Microbunches are produced by the interaction between the oscillating electrons in the undulator and the radiation that they produce (due to the oscillatory acceleration) leads to periodic longitudinal density modulation known as Microbunching. The in-phase emitted radiation adds coherently, increasing intensity and enhancing microbunching. Adapted from [1]. v

**Pulse** Also known as *Microbunch*. Each Train contains about 500 pulses produced from the self-amplified spontaneous emission (SASE) process (See Self-Amplified Spontaneous Emission). These are which are used as a secondary index for the data reduction process. 1, 19

**Self-Amplified Spontaneous Emission** is a process where the electron beam in the accelerator, when passing through an undulator, starts emitting radiation due to acceleration. The interaction between the emitted radiation and the charge distribution leads to microbunching. These microbunches emit radiation coherently, leading to the intense, coherent radiation, characteristic of a free-electron laser (FEL). For more information see [1]. v, 1

**Train** Also known as *Macrobunch*. A train represents a group of closely spaced electron bunches produced and accelerated by the FEL (or more generally, any accelerator). Each train is associated with a unique identifier called `trainId`, which is used as the primary index for much of the data reduction process. v, 1, 19

**Undulator** Magnets arranged periodically to produce a periodic magnetic field. It is used to produce coherent radiation by accelerating electrons through it. v,  
1

# Introduction

Conjecturing hypotheses based on observations has always been an important aspect of the scientific methodology. Especially in natural sciences, where the aim is to describe natural phenomena, the role of observations can not be overstated.

Trying to make sense of the observations is the.

After formulating a hypothesis, an experiment is designed to test it. From the evidence gathered through the experiment and through the aid of statistical principles, these hypotheses can be accepted (or rejected) with a defined level of confidence.

Hence, science is inherently linked to data. Data science helps us to formally handle this data, regardless of the domain. We want testable outcomes

With the dimensionality and size of data increasing, more sophisticated tools are necessary.

Especially in the field of neuroscience where experiments can not be performed easily or directly, a lot of statistical tools are employed.

In modern days, a paradigm shift has occurred where instead of trying to explicitly model the system we are trying to learn about, has also brought revolutions such as machine learning where we don't need to know the model itself and basically make the machine learn the non-linear model. This doesn't explain the process but allows us to perform for example inference. Learning is a sort of metamodel, where we don't need to know the theory itself but can still make predictions.

Expensive aspect! People don't use math methods first. Then we go into denoising

With the necessity to acquire data in so many dimensions, because we have low electron counts.

Free-Electron Laser Undulator Train Pulse Self-Amplified Spontaneous Emission



# Photoemission Spectroscopy

In the seminal paper by Einstein [3], that laid foundations to Quantum Mechanics, Einstein postulated that light is made of discrete quanta of energy  $E = h\nu$  to explain the observations by Hertz and J.J. Thompson; explaining the photoelectric effect. The effect can be described by the Equation 2.1 where  $E_e$  the emitted kinetic energy,  $h$  is the plank's constant,  $\nu$  the frequency of the incoming photon and  $\phi$  the material-specific work function (also known as binding energy).

$$E_e = h\nu - \phi \tag{2.1}$$

The equation describes how incident photons on a surface eject photoelectrons, provided the photon energy  $h\nu$  exceeds  $\phi$ . This also highlights that the emitted kinetic energy  $E_e$  does not depend on the photon flux (photon counts per second). However, the flux increases the total amount of electrons released from the material.

It is then apparent that the binding energy of electrons can be found by irradiating light onto the material and measuring the  $E_e$  of photoelectrons. photoemission/photoelectron spectroscopy (PES), is exactly such a technique that leverages this principle to probe the electronic structure of materials. While the above equation describes at what energies electron come out, it does not explain why

A complete treatment of photoemission process needs the quantum theory of light-matter interaction, and this shall be introduced where necessary.

## 2.1 Light-Matter Interaction

Where a complete description is not always necessary.

## 2.2 Light Sources

Table top, Synchrotrons, FELs. Time resolved PES and how it relates to acquisition times.

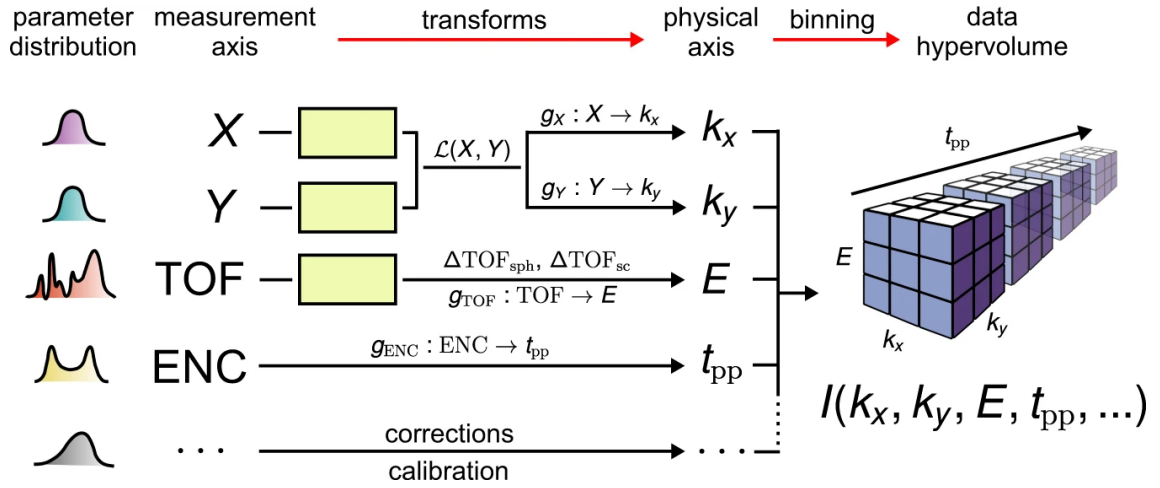


Figure 2.1: MPES taken from [8]

## 2.3 SASE FELs

## 2.4 HEXTOF Setup at FLASH

Deutsches Elektronen-Synchrotron (DESY) is a national

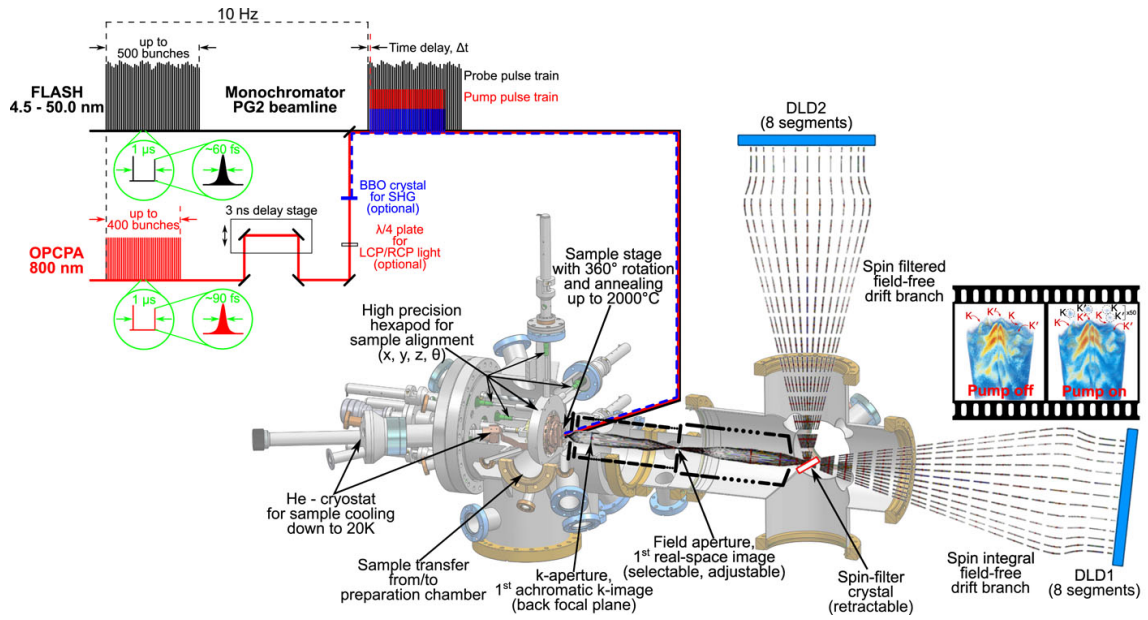


Figure 2.2: HEXTOF taken from [5]

# Imaging

## 3.1 Binning as means of Imaging

digital imaging from lecture sampling etc Binning is a way to find underlying distribution.

## 3.2 SASE fluctuations

Check the dld detector electron copies causing problem

## 3.3 Describing the data GdW, WSe<sub>2</sub>, GrIr, and new one

## 3.4 Creating the dataset Corrections Calibrations etc

# Characterizing Photon and Photo-electron Statistics

## 4.1 Poisson Process assumption

Poisson distribution is defined as [6]:

There's a lot of parameters that need to be tested to determine what sort of counting statistics the dataset has.

Controls for the test: lets see

- Total time being looked at (like 1000 s or 20 hours)
  - distribution might change due to overtime FEL intensity changes
- Time bins being used (like 2 s vs 20 s and so on).
  - Seems like distribution changes based on that too
- Looking at individual pixels on X and Y
- Looking at Energy axis as it behaves weirder
- Looking at a larger region in X and Y
  - should follow same statistics as single pixels
- Check after removing correlated electrons within each pulse
  - It is possible that electrons are correlated between different pulses because the time delay is long enough. But seems highly unlikely!

—

**Definition P1.20 (Poisson distribution; Siméon Poisson 1781-1840)**

The Poisson distribution  $\text{Poi}(\lambda)$  with parameter  $\lambda > 0$  is given by

$$P(\{k\}) = f(k) = \frac{\lambda^k}{k!} e^{-\lambda}, \quad k \in \mathbb{N}_0.$$

**Figure 4.1:** Enter Caption

Poisson data occur in all imaging processes where images are obtained by means of the count of particles, in general photons, arriving in the image domain. In the case of radioactive decay, fluorescence emission or similar phenomena, the arrival of particles is described by a *Poisson process*, i.e. a continuous stochastic process that is a collection of independent random variables  $\{N(t); t \geq 0\}$ , where  $N(t)$  is the number of particles arrived up to time  $t$ . The number of particles arriving within a given time interval  $T$  is a random variable (r.v.) with a *Poisson distribution* [52, 118], i.e. the probability of receiving  $n$  particles is given by

$$p(n) = \frac{e^{-\lambda} \lambda^n}{n!}, \quad n = 0, 1, 2, \dots, \quad (1)$$

where  $\lambda$ , proportional to  $T$ , is the expected value of the counts. This statistical model is appropriate to describe data acquired in fluorescence microscopy, emission tomography, optical/infrared astronomy, etc. Even if the energy of the photons (hence their wavelength) is different in the different applications, the statistics of the data is the same. The imaging system, however, significantly varies from an application to another. It is described by a sparse matrix in tomography, and by a convolution matrix in microscopy and astronomy.

Figure 4.2: Enter Caption

#### 4.1.1 Before filtering

#### 4.1.2 After filtering

Taken from [2] which cites [4]

## 4.2 Simulate Noise

We simulate the data with Poissonian data.

## 4.3 Statistical testing

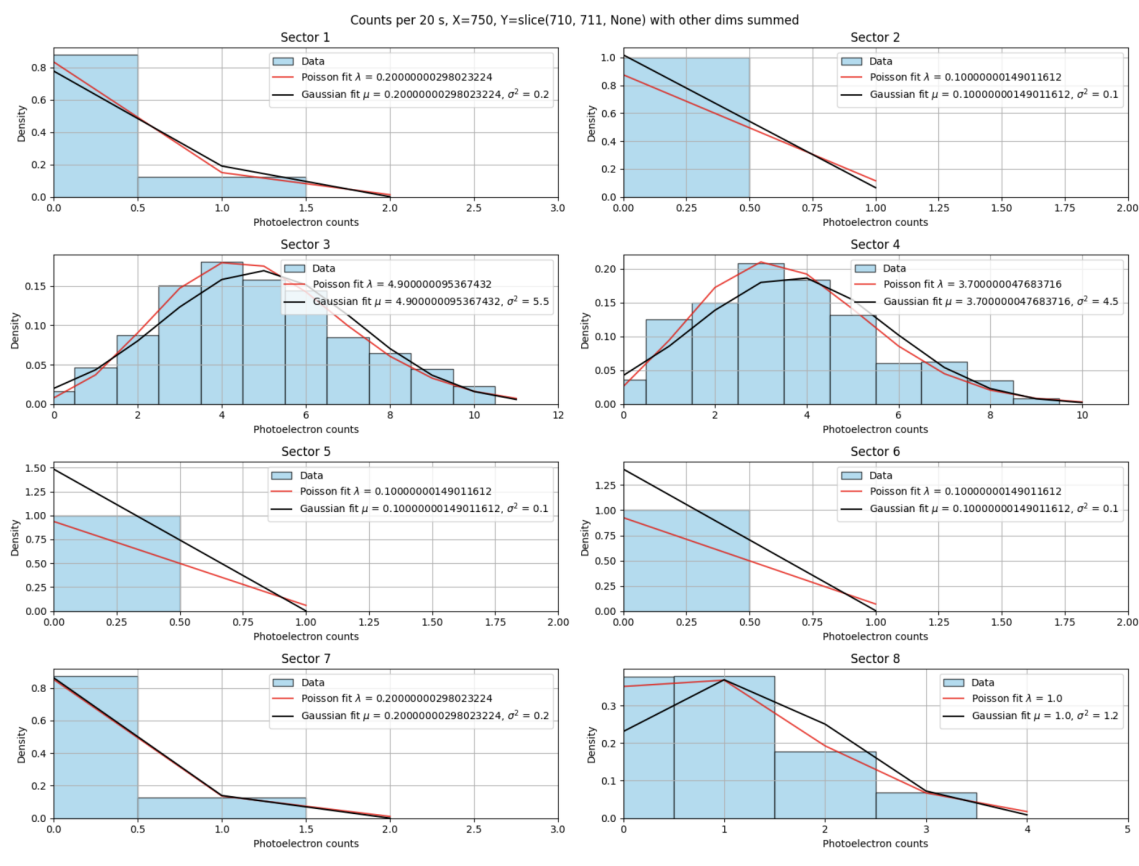
The **Central Limit Theorem** states that for a sequence of i.i.d. random variables  $X_1, X_2, \dots, X_n$  with mean  $\mu$  and variance  $\sigma^2$ , the normalized sum of these variables approaches a standard normal distribution as  $n$  tends to infinity:

$$\frac{\sum_{i=1}^n X_i - n\mu}{\sigma\sqrt{n}} \xrightarrow{d} N(0, 1) \quad \text{as } n \rightarrow \infty.$$

This convergence is in distribution.

From the law of large numbers, one can show that the relative fluctuations reduce as the reciprocal square root of the number of throws, a result valid for all statistical fluctuations, including shot noise. From Wikipedia

Shot noise exists because phenomena such as light and electric current consist of the movement of discrete (also called "quantized") 'packets'.



**Figure 4.3:** Enter Caption

### 4.3. Statistical testing

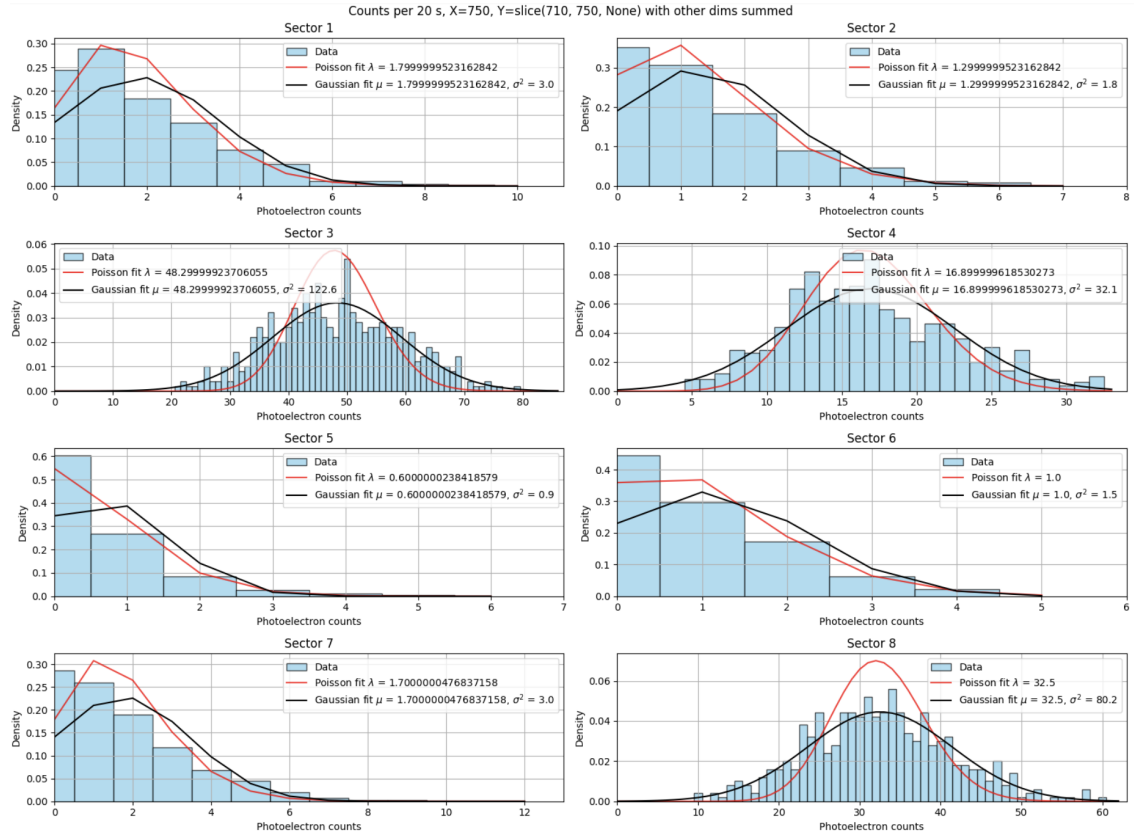


Figure 4.4: Image over 40 pixels and summed over energy

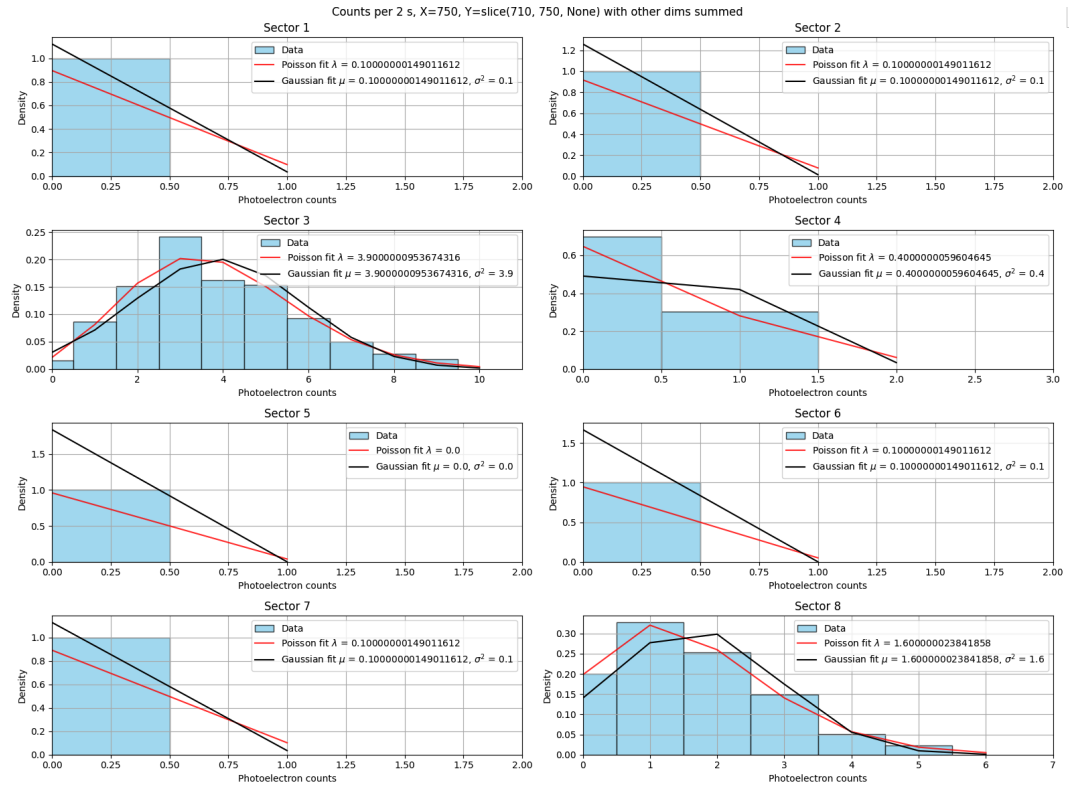
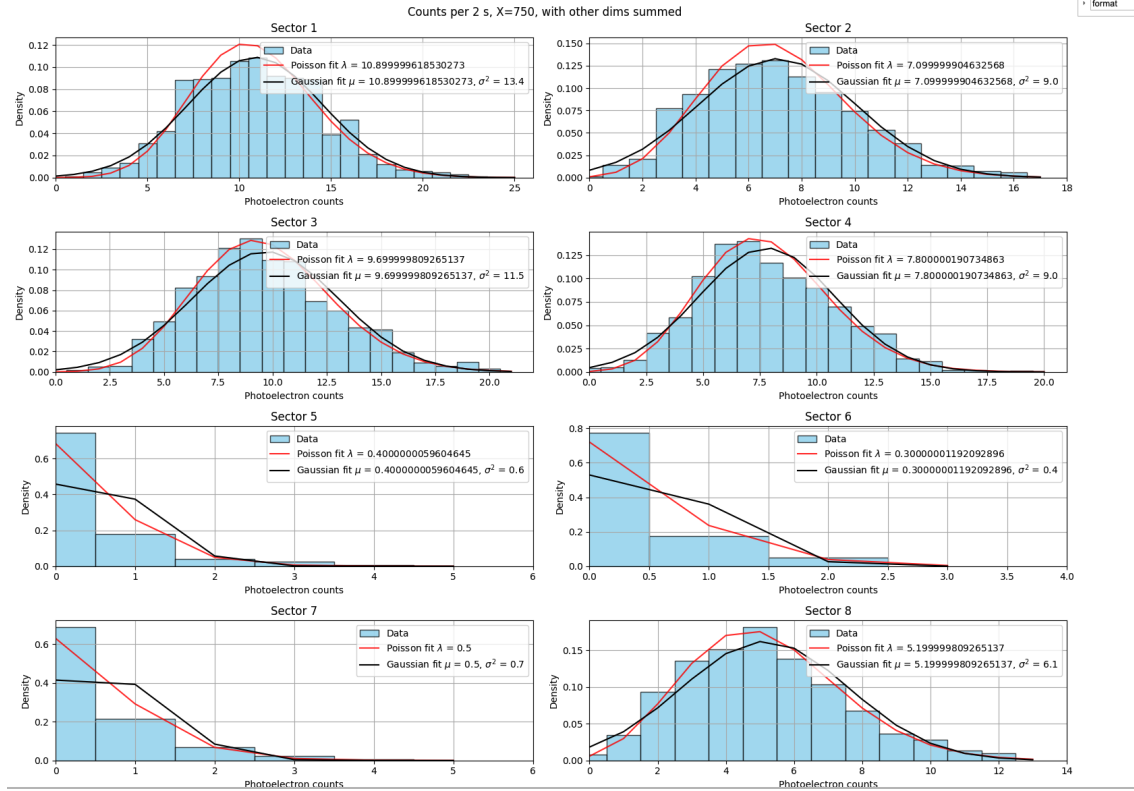


Figure 4.5: Data is filtered here with the KNN and looking at small region





**Figure 4.6:** Enter Caption

Can't correlate FEL intensity with electron counts per pulse as the GMD is before the monochromator.

The physical assumptions which we want to express mathematically are that the **conditions of the experiment remain constant** in time\*, and that non-overlapping time intervals are stochastically independent in the sense that information concerning the number of events in one interval reveals nothing about the other. The theory of probabilities in a continuum makes it possible to express these statements directly, but being restricted to discrete probabilities, we have to use an approximate finite model and pass to the limit.

Imagine a unit time interval partitioned into  $n$  subintervals of length  $1/n$ . A given collection of finitely many points in the interval may be regarded as the result of a chance process such that each subinterval has the same probability  $P_n$  to contain one or more points of the collection. A subinterval is then either occupied or empty, and the assumed independence of non-overlapping time intervals implies that we are dealing with Bernoulli trials: We assume that the probability for exactly  $k$  occupied subintervals is given by  $b(k; n, P_n)$ . We now refine this discrete model indefinitely by letting  $n \rightarrow \infty$ . The probability that the whole interval contains no point of the collection must tend to a finite limit. But this is the event that no cell is occupied, and its probability is  $(1 - p_n)^n$ . Passing to logarithms it is seen that this quantity approaches a limit

**Example S3.18 (Two-sample Kolmogorov-Smirnov test)**

Let  $X_1, \dots, X_{n_1} \stackrel{iid}{\sim} F$ , and  $Y_1, \dots, Y_{n_2} \stackrel{iid}{\sim} G$ , for strictly increasing, continuous CDFs  $F, G$ .

► Hypotheses:

$$H_0 : F = G \quad \leftrightarrow \quad H_1 : F \neq G$$

► Estimate  $F$  and  $G$  by

$$\begin{aligned}\hat{F}(x) &= \frac{1}{n_1} \sum_{i=1}^{n_1} \mathbb{1}_{(-\infty, x]}(X_i), & x \in \mathbb{R}, \\ \hat{G}(x) &= \frac{1}{n_2} \sum_{i=1}^{n_2} \mathbb{1}_{(-\infty, x]}(Y_i), & x \in \mathbb{R}.\end{aligned}$$

► Example P5.18:  $\|\hat{F} - F\|_\infty \xrightarrow{P} 0$  and  $\|\hat{G} - G\|_\infty \xrightarrow{P} 0$  as  $n_1, n_2 \rightarrow \infty$

► Test statistic:  $T_n = \|\hat{F} - \hat{G}\|_\infty$ .

**Figure 4.7:** Enter Caption

only if  $np_n$  from book [4]

## 4.4 Chi-squared Goodness of Fit Test

We hypothesize that the data follows a certain distribution e.g. Poisson, Normal, Negative Binomial. The Chi-squared Goodness of Fit Test is used to determine if the observed data is consistent with the expected distribution.

Let  $X_1, X_2, \dots, X_n \sim \text{i.i.d. } F$  and  $Y_1, Y_2, \dots, Y_m \sim \text{i.i.d. } G$ , where  $F$  and  $G$  are strictly increasing continuous CDFs.

The hypotheses for the chi-square test are:

$$H_0 : F = G \leftrightarrow H_1 : F \neq G \tag{4.1}$$

The hypothesis test for the Poisson distribution is the Chi-squared Goodness of Fit Test. The test statistic is given by:

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \tag{4.2}$$

where  $O_i$  is the observed frequency and  $E_i$  is the expected frequency. The degrees of freedom are given by  $k - 1$ , where  $k$  is the number of bins.

For goodness-of-fit tests, small p-values indicate that you can reject the null hypothesis and conclude that your data were not drawn from a population with the specified distribution. Consequently, goodness-of-fit tests are a rare case where you look for high p-values to identify candidate distributions. (From This website)

# Image Reconstruction

Rather than calling it denoising, better word in image reconstruction because Image Reconstruction:

Purpose: To reconstruct an image from incomplete, noisy, or indirect measurements. This is often used in medical imaging (e.g., MRI, CT scans), computational photography, and computer vision applications.

Reconstruction involves generating a complete image from partial or indirect data, which can include denoising and deblurring as sub-tasks.

## 5.1 Address reconstruction/denoising schemes

VST with BM3D: BM3D uses collaborative filtering, which is also used in recommender systems [citation need] UNET noise2noise

## 5.2 BM3D: Denoising in sparse domain

with and without anscombe on different datasets from flash, lab and fhi Testing s

### 5.2.1 Anscombe: Variance Stabilization Transform

## 5.3 Statistical Learning

One can do a lot of different learning approaches. We will build towards the deep learning approach. Networks which can deal with image data are generally Convolutional Neural Networks. Looking at research, UNET has been used for image segmentation and denoising, which combines the concept of autoencoders and convolutional neural networks, along with skip connections. We use this approach with the noise2noise framework, which is a deep learning framework for image reconstruction.

ERM Deep learning is just linear separator problem with a non linear function applied to it like RELU

Our aim is to reduce the error between the true image and the reconstructed

image. This is a regression problem, where we are trying to learn a function that maps the noisy image to the true image. Mathematically, this can be written as:

ERM: Empirical Risk Minimization

$$\hat{f} = \arg \min_f \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 \quad (5.1)$$

Popular loss functions include the mean squared error (MSE), mean absolute error (MAE), and Huber loss. The choice of loss function depends on the noise model and the desired properties of the reconstruction. For example, the MSE is commonly used for Gaussian noise, while the MAE is more robust to outliers.

The MSE is defined as:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 \quad (5.2)$$

MSE is the squared L2 norm of the difference between the predicted and true values. It is sensitive to outliers and can be dominated by large errors. The MAE is defined as:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |f(x_i) - y_i| \quad (5.3)$$

### 5.3.1 Optimal Loss Function

Optimal loss can be derived from the distribution model. The Method of Moments and Maximum Likelihood Estimation are two common methods for deriving the optimal loss function.

For Poisson noise, the optimal loss function is the negative log-likelihood of the Poisson distribution. This is because the Poisson distribution is the maximum entropy distribution for count data, and the negative log-likelihood is the maximum likelihood estimator for the Poisson distribution. This can be written as the following optimization problem:

$$\hat{f} = \arg \min_f - \sum_{i=1}^n \log \left( \frac{e^{-f(x_i)} f(x_i)^{y_i}}{y_i!} \right) \quad (5.4)$$

### 5.3.2 Regularization

There are different ways to regularize the loss function to prevent overfitting. This can be done by adding a penalty term to the loss function. Common regularization techniques include L1 and L2 regularization, which add the absolute value of the weights and the square of the weights to the loss function, respectively. This can be written as:

$$\hat{f} = \arg \min_f \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 + \lambda \sum_{j=1}^p |\beta_j| \quad (5.5)$$

### **5.3.3 Optimization**

## **5.4 Noise2Noise: Deep Learning framework**

Important to take care not to train with empty data.

### **5.4.1 Convolutional Neural Networks**

### **5.4.2 Autoencoder**

# Conclusion and Outlook

We use SASE currently but with seeded FEL, might be more robust poisson statistics. Maybe even sub-poissonian.

# Bibliography

- [1] W. Ackermann et al. “Operation of a Free-Electron Laser from the Extreme Ultraviolet to the Water Window”. In: *Nature Photonics* 1.6 (June 2007), pp. 336–342. ISSN: 1749-4893. DOI: 10.1038/nphoton.2007.76. (Visited on 08/25/2024).
- [2] M Bertero et al. “Image Deblurring with Poisson Data: From Cells to Galaxies”. In: *Inverse Problems* 25.12 (Dec. 2009), p. 123006. ISSN: 0266-5611, 1361-6420. DOI: 10.1088/0266-5611/25/12/123006. (Visited on 08/19/2024).
- [3] A. Einstein. “Über Einen Die Erzeugung Und Verwandlung Des Lichtes Betreffenden Heuristischen Gesichtspunkt”. In: *Annalen der Physik* 322.6 (1905), pp. 132–148. DOI: 10.1002/andp.19053220607.
- [4] William Feller. *An Introduction to Probability Theory and Its Applications*. Third ed. rev. Wiley Series in Probability and Mathematical Statistics. New York Chichester Brisbane [etc.]: J. Wiley, 1968. ISBN: 978-0-471-25708-0.
- [5] D. Kutnyakhov et al. “Time- and Momentum-Resolved Photoemission Studies Using Time-of-Flight Momentum Microscopy at a Free-Electron Laser”. In: *Review of Scientific Instruments* 91.1 (Jan. 2020), p. 013109. ISSN: 0034-6748, 1089-7623. DOI: 10.1063/1.5118777. (Visited on 08/19/2024).
- [6] Siméon-Denis Poisson. *Recherches sur la probabilité des jugements en matière criminelle et en matière civile: précédées des règles générales du calcul des probabilités*. Bachelier, 1837.
- [7] Muhammad Zain Sohail. “Ultrafast Dynamic Studies in Spin-Filter Material Europium Oxide with Resonant Inelastic X-ray Scattering”. Bachelor’s Thesis. Carl von Ossietzky Universität Oldenburg, 2021.
- [8] R. Patrick Xian et al. “An Open-Source, End-to-End Workflow for Multidimensional Photoemission Spectroscopy”. In: *Scientific Data* 7.1 (Dec. 2020), p. 442. ISSN: 2052-4463. DOI: 10.1038/s41597-020-00769-8. (Visited on 08/25/2024).

# Acknowledgements

Thank Dima especially Dataset providers Thank Lorenz Kruger, Martin and others in group (Lukas Weigand). This research was supported in part through the Maxwell computational resources operated at Deutsches Elektronen-Synchrotron DESY, Hamburg, Germany



# Contributions

A preliminary analysis during the early works of this thesis titled *Efficient Data Acquisition in Multi-Dimensional Photoemission Spectroscopy using Denoising*, was presented, in the form a of poster, at the DPG Spring Meetings 2024, as well as NanoMat Science Day 2024 at DESY.

# Declaration

Insert declaration here Date & Signature: \_\_\_\_\_

# Appendices

# Transforming Raw Data to structured format

## Extract, Transform, Load

Raw data from the experiment is stored in hierarchical data format version 5 (HDF5) files. This includes many Beamline diagnostic information such as beam arrival monitor (BAM), gas monitor detector (GMD), the delay stage readings, the monochromator energy, sample specific information such as extractor voltage, and the three-dimensional electron counting by the delay line detector (DLD). This information is resolved at each bunch of electrons coming from the accelerator called a Train, which are further microbunched into pulses.

The open community of multidimensional photoemission spectroscopy (OpenCOMPES) was established to develop tools and infrastructure to make analysis easier. To this end, a modular Python library called **Single Event DataFrame (SED)** was created that provides the entire pipeline from easy data ingestion to common calibration and corrections, Multidimensional binning to create images, and saving the images standard formats, with data provenance.

The data pipeline is designed to extract raw data from HDF5 files, transform the data into a structured format for analysis, and load the transformed data into buffer files. These buffer files are subsequently used for downstream processing, analysis, and visualization. The following sections provide a detailed description of each stage in the ETL process.

## Pipeline Overview

The ETL pipeline is divided into several key stages, each critical to the preparation and validation of the data. These stages include:

- **Data Extraction:** Retrieving raw H5 files from experimental runs.
- **Buffer File Creation:** Generating interim buffer files that facilitate further processing.

- **Data Transformation and Validation:** Applying domain-specific transformations, such as forward-filling non-electron channels and splitting sector IDs, while validating data integrity against predefined schemas.
- **Data Structuring:** Organizing the processed data into structured Parquet files suitable for downstream analysis.

## Detailed Pipeline Stages

### Data Extraction

The initial stage of the pipeline involves the extraction of raw data from H5 files. These files contain experimental data with various channels, such as electron and timed information, stored in a hierarchical structure. The paths to these H5 files are provided as input to the pipeline, along with configuration settings that dictate the subsequent processing steps.

### Buffer File Creation

To streamline the data processing, the pipeline first creates buffer files for each type of data (electron and timed dataframes). This is achieved through the `BufferFilePaths` class, which initializes paths for the raw H5 files and corresponding buffer files. The class checks for existing buffer files and determines whether new files need to be generated or existing ones should be reused, based on the `force_recreate` flag.

For each H5 file, the pipeline generates two buffer files:

- **Electron Buffer File:** Contains data relevant to individual electron events.
- **Timed Buffer File:** Aggregates data at pulse and train levels, resolving timing information without individual electron data.

### Data Transformation and Validation

Once the buffer files are established, the pipeline proceeds to transform the data. The key transformation steps include:

- **Forward-Filling Non-Electron Channels:** Missing values in non-electron channels are filled using a forward-fill strategy to ensure data continuity.
- **Schema Validation:** The schema of the generated Parquet files is validated against the expected schema derived from configuration files. This ensures that all required channels are present and correctly formatted. The schema check involves:
  - Reading the schema from existing Parquet files.
  - Comparing the actual schema with the expected schema.

- 
- Raising errors if discrepancies are found, prompting a review of the configuration or a forced recreation of buffer files.
  - **Splitting Sector ID from DLD Time:** A custom transformation is applied to separate the sector ID from the DLD (Delay Line Detector) time within the electron dataframe. This operation is essential for accurately resolving electron events.

## Data Structuring and Finalization

After transforming and validating the data, the pipeline structures it into Parquet files, which are columnar storage formats optimized for analytical queries. The steps involved include:

- **Electron Dataframe Structuring:** The electron-resolved dataframe is processed to drop non-electron data and reset the index. The processed data is then saved as a Parquet file.
- **Timed Dataframe Structuring:** The timed dataframe is derived by aggregating data at pulse and train levels, excluding electron-specific data. This structured data is also saved as a Parquet file.
- **Metadata Generation:** Metadata related to file statistics, filling operations, and schema checks is compiled and saved, providing crucial information for downstream analysis and data auditing.

# Mathematical Background

## B.1 Measure Space and Measures

### Measurable Space.

Let  $\Omega = \emptyset$ ,  $P(\Omega)$  the power set of  $\Omega$  and  $\mathcal{A} \subset P(\Omega)$ . If the following conditions are met,  $\mathcal{A}$  is called  $\sigma$ -algebra, and the  $(\Omega, \mathcal{A})$  pair make up a measurable space.

1.  $\Omega \in \mathcal{A}$ .
2. If  $A \in \mathcal{A}$ , then  $\Omega \setminus A \in \mathcal{A}$ .
3. If  $(A_n)_{n \in \mathbb{N}}$  is a sequence of sets where  $A_n \in \mathcal{A}$  for all  $n \in \mathbb{N}$ , then

$$\bigcup_{n \in \mathbb{N}} A_n \in \mathcal{A}.$$

**Positive Measure.** Let  $(\Omega, \mathcal{A})$  be a measurable space as defined in above. A function  $\mu : \mathcal{A} \rightarrow [0, \infty]$  is called a *positive measure* if it satisfies the following conditions:

1.  $\mu(\emptyset) = 0$ , (the measure of the empty set is zero)
2.  $\mu$  is *countably additive*: For any countable collection of disjoint sets  $(A_i)_{i \in \mathbb{N}} \subset \mathcal{A}$ , we have

$$\mu \left( \bigcup_{i=1}^{\infty} A_i \right) = \sum_{i=1}^{\infty} \mu(A_i).$$

If these conditions are met, then  $\mu$  is called a *measure* on the measurable space  $(\Omega, \mathcal{A})$ .

## B.2 Probability

The theory of probability is necessary to quantify stochastic and uncertain quantities. Throughout this text, it can be seen used in the context of inherently stochastic processes, such as quantum effects, and to quantify uncertainty in measurements.

**Probability Measure.** A *probability measure*  $P : \mathcal{A} \rightarrow [0, 1]$  satisfies all the properties of a positive measure (see Section B.1), with the additional property known as the normalization condition:

$$P(\Omega) = 1$$

This leads to the measure space for probability (or probability space) being  $(\Omega, \mathcal{A}, P)$ . The probability of an event  $A \in \mathcal{A}$  is  $P(A)$ . The probability measure is hence a real number between 0 and 1 that defines the likelihood of an event to occur.

In the *frequentist* interpretation, probability is defined as the long-run relative frequency of an event occurring in repeated independent trials. It assumes that probabilities are objective and intrinsic properties of the physical world. For example, the probability of getting heads in a fair coin toss is 0.5, meaning that if we were to toss the coin an infinite number of times, half of the outcomes would be heads.

The *Bayesian* interpretation, on the other hand, views probability as a measure of belief or certainty about an event, given the available information. It is inherently subjective and updates as new evidence is introduced. For instance, if we initially believe that a coin is fair, but after observing several tosses we notice a bias, we update our belief (and hence the probability) to reflect the new evidence. This process of updating beliefs is formalized through *Bayes' rule*, which is a cornerstone of Bayesian inference.

**Law of Total Probability.** The law of total probability provides a way to compute the probability of an event based on a partition of the sample space. If  $\{B_i\}_{i=1}^n$  is a partition of the sample space  $\Omega$  (i.e.,  $B_i \cap B_j = \emptyset$  for  $i \neq j$  and  $\bigcup_{i=1}^n B_i = \Omega$ ), then for any event  $A$ :

$$P(A) = \sum_{i=1}^n P(A \mid B_i)P(B_i). \quad (\text{B.1})$$

This law is particularly useful when dealing with complex events that can be decomposed into simpler, mutually exclusive cases.

## B.3 Landau Notation

Landau notation, commonly referred to as Big O notation and its relatives (Big Omega, Big Theta, etc.), is used to describe the asymptotic behavior of functions.



This is particularly important in the analysis of algorithms and in expressing the growth rates of functions.

## **Big O Notation ( $\mathcal{O}$ )**

### **B.4 Statistical Inference**

#### **B.4.1 Testing/Confidence Intervals**

Testing if distribution is poissonian

#### **B.4.2 Optimal parameter estimation**

Used in estimating distribution parameters

# Deep Learning

## C.1 Infrastructure

Model Architecture: UNET 2D and 3d

Trained on V100 GPU using Maxwell Cluster at DESY.