

MASTER OF ENGINEERING PROJECT REPORT

Characterization of LNP Drug Delivery Vehicles by Machine Learning SAXS Part

Student:

Tianyou Li (NetID: tl899)

Project Advisor:

Prof. Peter Doerschuk

Department of Electrical Computer Engineering

Contents

1	Introduction	4
2	Methodology	4
2.1	Data Generation	4
2.1.1	3D Model Preparation	5
2.1.2	Adaptive Padding + FFT	7
2.1.3	Continuous Density Correction	10
2.1.4	Orientation Averaging + Spherical Correction .	12
2.1.5	Data Pre-processing	17
2.1.6	Data Synthesizing	19
2.1.7	Within-Type Data Augmentation	19
2.1.8	Cross-Type Data Synthesizing	20
2.1.9	Poisson Noise Addition	22
2.2	ML Model Design	23
3	Results	25
3.1	Results for Classification	25
3.2	Results for Fraction Regression	27
4	Conclusion	27

Abstract

This study builds upon existing research by leveraging machine learning (ML) techniques to enhance the characterization of lipid nanoparticles (LNPs) for drug delivery applications. By integrating physics-based simulation methods, advanced experimental data preprocessing, and realistic noise modeling, we have developed highly accurate small-angle X-ray scattering (SAXS) simulations for dilute aqueous LNP solutions. Our approach employs convolutional neural networks (CNNs) and residual neural networks (ResNets) to effectively classify heterogeneous LNP solutions and precisely predict size distribution parameters for homogeneous systems. This methodology offers a rapid and cost-effective alternative to conventional techniques such as cryo-electron microscopy (cryo-EM), thereby holding significant promise for advancing nanomedicine.

Keywords: Lipid Nanoparticles, SAXS, Physics-based Simulation, Machine learning

1 Introduction

Nanoparticle-based drug delivery systems—particularly lipid nanoparticles (LNPs)—have revolutionized therapeutic delivery by enabling highly precise and efficient treatment modalities. LNPs garnered significant global attention during the COVID-19 pandemic due to their critical role in mRNA vaccine development. Their demonstrated versatility, stability, and capacity for targeted tissue delivery have established them as pivotal components in advancing precision medicine.

Accurate characterization of the nanoscale structural and functional properties of LNPs is essential for optimizing their therapeutic potential. Two techniques are commonly employed for this purpose: cryo-electron microscopy (cryo-EM) and small-angle X-ray scattering (SAXS). Cryo-EM provides high resolution, two-dimensional projections of individual LNPs; however, its high operational costs, slow throughput, and limited integration into production environments constrain its routine application. In contrast, SAXS yields an ensemble-averaged one-dimensional scattering profile, which—although offering less direct insight into individual particle morphology—is relatively inexpensive, rapid, and well-suited for incorporation into production lines. Therefore, the objective of this work is to assess the extent to which ML can extract detailed structural information from SAXS data, using cryo-EM as the benchmark standard.

2 Methodology

2.1 Data Generation

This subsection provides a comprehensive overview of the data generation pipeline, as depicted in Fig. 1. The following sections detail each step of the process. In brief, our methodology encompasses the preparation of 3D LNP models, the implementation of an adaptive padding strategy combined with FFT, and the application of a continuous density correction to produce realistic SAXS simulations. Each stage is carefully designed to capture the relevant physical phenomena while ensuring computational efficiency and fidelity to experimental conditions.

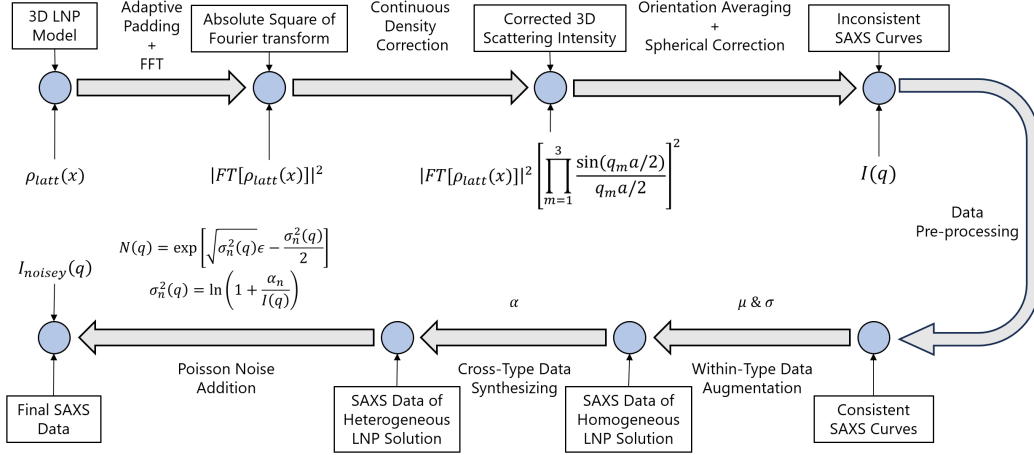


Figure 1: The Data Generation Pipeline Flow Chart. Symbols: $\rho_{latt}(x)$ denotes the 3D LNP model; FT denotes the results of the Fourier Transform; q_m denotes the minimum q value in the reciprocal space; a denotes the real space resolution in nm; μ, σ denote the random variables defining a unique lognormal distribution; α denotes the mixing factor (fraction of Type 3 LNPs); σ_n denotes the q -dependent variance; ϵ denotes a standard normally distributed random variable; α_n denotes a log-uniform scaling factor.

2.1.1 3D Model Preparation

Before diving into the details of how to generate realistic SAXS data, it is also necessary to ensure the 3D LNP models adopted in this project are valid, which means the shapes or sizes should be verified by cryo-EM data. As shown in Fig. 2, it is not suitable to assume the LNPs are spherical, which, however, is the normal practice when analyzing the SAXS data. Therefore, in this project, we try to apply more detailed 3D LNP models to generate the SAXS data and see whether the ML techniques can extract structural information from the generated SAXS data.

The LNP models are constructed to mirror the actual physical structure observed in cryo-EM images and SAXS experiments. Specifically, the models adopt a core-shell ellipsoid design where the core is defined by an equatorial radius that ranges between 21 and 30 nm and an axial ratio of about 1.667—parameters determined from detailed image measurements that shifted the model from an initially assumed oblate shape to a prolate one. The shell, representing the lipid layers, is given a uniform thickness of approximately 4 nm, an intermediate value chosen from experimental estimates between 3 and 5 nm. In MATLAB simulations, this physical structure is further detailed by mapping the electron density in a 3D grid (using a density function that is

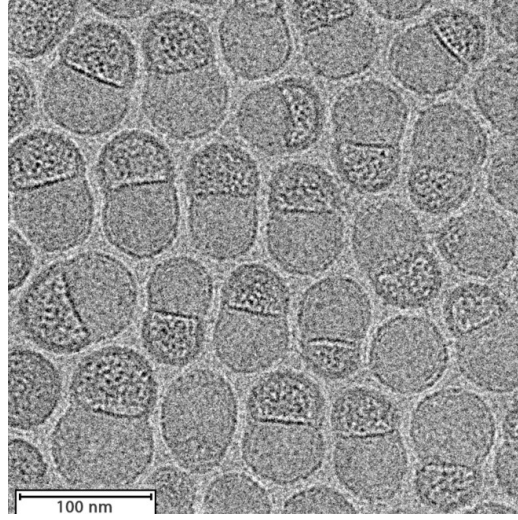


Figure 2: Cryo-EM image of mRNA-LNP

constant inside the particle and zero outside) and performing Fourier transforms to generate SAXS curves. To further diversify the dataset, a scaling factor is employed to control the size of 3D models, resulting in the simulated LNPs' sizes ranging from 20 to 100 nm .

At this stage of the project, we created 4 types of different models as shown in Fig. 3, both of which can be found in either Fig. 2 or other cryo-EM data sources. Under the current settings, the model is $100 \times 100 \times 100$ in total, where each voxel corresponds to 1 nm^3 cube in real space. According to [1], LNPs in the size range of 40 to 80 nm are more desired in production. Therefore, both the classifier and the prediction model are designed to classify the solution type and predict the size distribution within the range of 20 to 100 nm . Thus, there are 81 models in total with 1 nm increment in diameter, which are regarded as ground truth in the following data synthesizing process. These models will be used for SAXS data simulation since they represent the electron density of the LNPs. Hence, we have obtained:

$$\rho_{latt}(x)$$

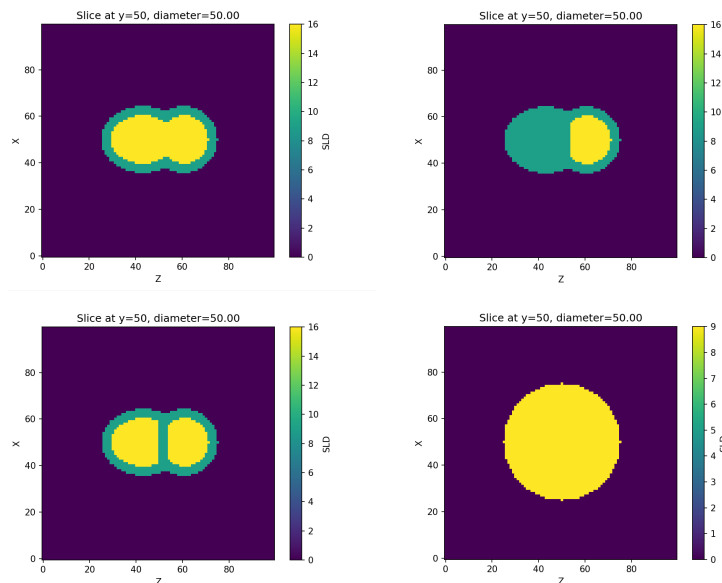


Figure 3: Four Different Types of 3D LNP Mmodel. SLD denotes the scattering length density, reflecting the material’s scattering ability to X-rays or neutron beams

2.1.2 Adaptive Padding + FFT

In SAXS simulations using a finite computational cube, discontinuities at the boundaries can introduce artifacts in the Fast Fourier Transform (FFT). These discontinuities result in distortions and spurious oscillations in the computed scattering intensity. The FFT inherently assumes periodic input data, effectively replicating the electron density cube in all directions (as illustrated in Fig. 3). When the density at the boundaries does not seamlessly match the surrounding values (which, for an isolated particle, is zero), a sharp discontinuity occurs as the data “wraps around.” Such abrupt transitions, akin to step functions in real space, inherently contain a wide range of frequencies—a phenomenon related to the Gibbs effect. Consequently, power “leaks” into many frequencies, manifesting as spectral leakage that appears as noise or artificial ripples in the output. For example, Fig. 4 compares simulated SAXS data for a 50 nm LNP model using two different padding sizes (200 nm and 600 nm), clearly demonstrating that reduced padding exacerbates distortion.

Although windowing techniques can smooth the edges to mitigate these effects, they unavoidably alter the amplitude and finer details of the signal.

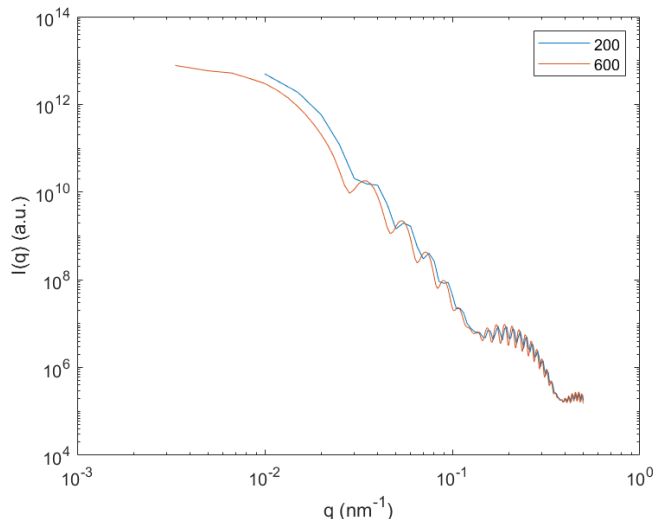


Figure 4: Illustration of the FFT Boundary Effect

To preserve the physical accuracy of the simulation, the optimal approach is to enlarge the computational domain so that the electron density naturally decays to zero at the boundaries, thereby avoiding the need for artificial windowing. A common strategy for simulating an isolated particle is zero-padding, where the cube is extended with zeros (representing a vacuum) around the particle. This method increases the domain size and separates the particle from its periodic images, ensuring a smooth transition to zero at the boundaries. A widely adopted guideline is to restrict the particle to no more than half the linear dimension of the computational box to prevent overlap with its periodic replicas.

In this project, the primary goal is to design a ML algorithm that extracts genuine physical information from the simulated SAXS data—such as distinguishing between LNP solution types or predicting the mean size—without relying on numerical artifacts. Employing a fixed padding size would result in different relative padding for small and large LNPs, yielding FFTs with different absolute resolutions. Such variations could inadvertently provide cues to the ML model that are unrelated to the intrinsic properties of the LNPs.

To overcome this challenge, we adopt an adaptive padding strategy that

maintains a fixed ratio between the padded cube size and the effective electron density cube. This approach offers several benefits:

- **Uniform Relative Boundary Effects:** A constant padding-to-model ratio ensures that the effective distance from the model to the boundaries is proportional for all samples, leading to consistent truncation-induced discontinuities and spectral leakage.
- **Consistent Relative FFT Resolution:** Although the absolute FFT resolution (i.e., the spacing in reciprocal space) varies with the overall cube size, adaptive padding ensures that the resolution remains consistent relative to the model dimensions. This consistency is preserved when interpolating to a unified reciprocal axis, thereby maintaining the physical features of the scattering data.
- **Prevention of Spurious Cues for ML:** By standardizing the relative padding, the FFT artifacts (e.g., due to boundary discontinuities) are uniformly distributed, reducing the risk that the ML algorithm will associate these numerical differences with intrinsic particle properties.
- **Physical Consistency:** Scaling the simulation environment in proportion to the effective model better represents the natural decay of electron density, preserving the physical integrity of the simulation.

Under this strategy, for example, a 50 nm electron density cube is padded to 300 nm, and a 100 nm cube to 600 nm before performing the FFT. This uniformity in relative boundary effects compels the ML model to rely on features that genuinely reflect the LNP models' shape and size, rather than on numerical artifacts arising from differences in FFT resolution.

Before implementing the correction algorithm, it is necessary to establish the method for constructing reciprocal (q) space. In our approach, the padding size is specified in terms of grid points rather than physical units. For instance, if the original electron density cube consists of 50 points, a padded cube may be defined with 300 points (denoted as n_x) to capture maximal information from the real-space model. The minimum and maximum q values are then defined as:

$$q_{\min} = q^{(1)} = \frac{2\pi}{n_x a}, \quad q_{\max} = q^{N/2} = \frac{\pi}{a}, \quad \text{where } a = 1 \text{ nm.}$$

The maximum q is defined as $q^{N/2}$ because the FFT output is shifted so that the zero-frequency component is centered—a necessary condition for orientation averaging. With the padded cube prepared, a three-dimensional FFT is performed. Notably, the number of FFT points is not constrained to be a power of 2, as empirical tests have shown negligible performance loss. The resulting quantity is given by:

$$|FT[\rho(x)]|^2.$$

2.1.3 Continuous Density Correction

In numerical simulations, the electron density is defined only at discrete lattice points, whereas the physical electron density is continuous. To generate the most realistic SAXS data, it is therefore necessary to implement a continuous density correction [2]. As illustrated in Fig. 5, this correction bridges the gap between discrete density peaks and a continuous density distribution by convolving the discrete peaks (as used in numerical DFT) with an elementary “box” function of width a [2].

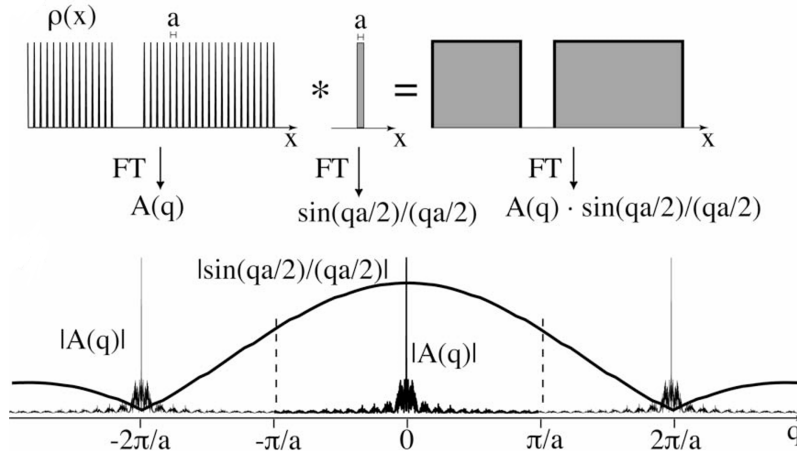


Figure 5: Continuous-Discrete Density Distribution Correction Symbols: $A(q)$ denotes the result of Discrete Fourier Transform [2]

Mathematically, the Fourier Transform (FT) of a convolution is equal to the product of the FTs of the individual functions. Thus, the corrected scattering

intensity can be expressed as:

$$I(\mathbf{q}) = |FT[\rho_{\text{latt}}]|^2 \left[\prod_{m=1}^3 \frac{\sin(q_m a/2)}{q_m a/2} \right]^2 \quad [2]$$

where ρ_{latt} denotes the discretely defined electron density. In other words, we approximate the scattering intensity obtained from the continuous density by multiplying the intensity from the discrete density by the correction factor

$$\left[\frac{\sin(q m_a/2)}{q m_a/2} \right]^2$$

in each dimension ($m = 1, 2, 3$). This procedure is illustrated in Fig. 5 for a one-dimensional case. For discrete q values, we define

$$q^{(k)} = \frac{2k}{Na},$$

and multiplying by $a/2$ yields

$$\frac{q^{(k)} a}{2} = \frac{k}{N}$$

Thus, along one dimension, the k th point (out of N) in the IDFT(q) data field is multiplied by

$$\left\{ \frac{\sin \left[\frac{(k-N/2-1)}{N} \right]}{\frac{(k-N/2-1)}{N}} \right\}^2$$

taking into account that the point corresponding to $q = 0$ is located at $N/2 + 1$. In three dimensions, with a lattice spacing given by

$$q_s = \frac{2}{Na}$$

the continuous scattering intensity is obtained via

$$I_{\text{cont}}(n_1 q_s, n_2 q_s, n_3 q_s) = I_{\text{DFT}}(n_1 q_s, n_2 q_s, n_3 q_s) \prod_{m=1}^3 \left\{ \frac{\sin \left[\frac{(n_m-N/2-1)}{N} \right]}{\frac{(n_m-N/2-1)}{N}} \right\}^2 \quad (1)$$

For small q (or equivalently, for large structures, which are most relevant in small-angle scattering), the sinc functions in Equation (1) are nearly unity and leave $I(q)$ virtually unchanged. This is expected since small q values correspond to large distances where the fine details of the electron density do not significantly affect the scattering intensity. At the limit of the first Nyquist zone, $|q| = \pi/a$, the scattering intensity $I(\pi/a)$ is reduced by a factor of $(2/\pi)^2 \approx 0.4$. Therefore, the corrected scattering intensity is given by:

$$|FT[\rho_{latt}]|^2 \left[\prod_{m=1}^3 \frac{\sin(q_m a/2)}{q_m a/2} \right]^2$$

2.1.4 Orientation Averaging + Spherical Correction

After applying the necessary correction functions, the next step is to perform orientation averaging to convert the full three-dimensional scattering data into a one-dimensional intensity curve, $I(q)$. The fundamental idea is to compute the scattering intensity at a given q by averaging the intensities of all points in 3D q space with $|\mathbf{q}| \approx q$. This is equivalent to averaging over a spherical shell in q space, hence the term “Orientation Averaging” or “Spherical Mean”. In our implementation, the FFT yields n_x^3 points in the 3D q space, after which a one-dimensional q axis with n_{xf} uniformly spaced bins is defined. The following discussion first outlines the original algorithm proposed in [2] and then introduces an improved algorithm along with a comparative analysis.

Original Algorithm

The primary objective of the original algorithm is to distribute the scattering intensity from each of the n_x^3 points onto the n_{xf} bins. Since the bins are uniformly spaced over the overall q range, each value $q = |\mathbf{q}|$ must be mapped onto this discretized axis. The process can be summarized as follows. First, the n_x^3 points are flattened into a one-dimensional array, and the q value for the κ^{th} point is computed as:

$$q = |\mathbf{q}| = \frac{2\pi\kappa}{n_x a},$$

where a denotes the unit length. With the new 1D q axis consisting of n_{xf}

bins, let k be the bin index corresponding to the nearest value

$$q' = \frac{2\pi k}{n_{xf} a}$$

to $\frac{2\pi\kappa}{n_x a}$. Because q generally lies between two consecutive q' values, the intensity $I_{\text{DFT}}(\mathbf{q})$ is apportioned between the adjacent bins. Specifically, a fraction

$$1 - \frac{|q' - q|}{\Delta q} = 1 - |k - \kappa|, \quad \text{where } \Delta q = \frac{2\pi}{n_{xf} a},$$

is allocated to the bin at

$$q' = \frac{2\pi k}{n_{xf} a},$$

while the remaining fraction

$$\frac{|q' - q|}{\Delta q} = |k - \kappa|$$

is assigned to the adjacent bin at

$$q' = \frac{2\pi (k + \text{sgn}(q' - q))}{n_{xf} a}.$$

In formal terms, this process can be written as:

$$(1 - |\kappa - k|) I_{\text{DFT}}(\mathbf{q}) \rightarrow I' \left(\frac{2\pi k}{n_{xf} a} \right), \quad |\kappa - k| I_{\text{DFT}}(\mathbf{q}) \rightarrow I' \left(\frac{2\pi (k + \text{sgn}(\kappa - k))}{n_{xf} a} \right).$$

Since the FFT of real-valued functions is symmetric about zero frequency, this “channel sharing” process is performed for every point in one-eighth of the n_x^3 grid, effectively accumulating their contributions into the n_{xf} uniformly distributed bins. Additionally, to compensate for the uneven number of points in different regions, the central point is assigned a weight of 1 while other points are weighted by 2, reflecting their double counting in the complete n_x^3 dataset.

Because the number of points with $q = |\mathbf{q}|$ in a spherical shell increases proportionally to the square of the radius (as noted in [2] and [3]), the aggregated

scattering intensity from the cubic lattice is given by:

$$I'(q) = I(q)q^2$$

Without proper correction, the scattering intensities at higher q values will be overestimated due to the larger number of points in the spherical shell, leading to a statistical bias. Therefore, to obtain the correct scattering intensity, $I'(\mathbf{q})$ must be divided by q^2 to yield $I(\mathbf{q})$ [2].

Fig. 6 compares SAXS data generated for a spherical particle using the original algorithm with an analytical solution provided by Professor Peter Dorschuk. The analytical result, derived from a simple mathematical model, serves as a benchmark since spherical particles are the only case for which an analytical expression is available. The comparison indicates that while the original algorithm approximates the scattering intensity well at high q , discrepancies exist in the low q region.

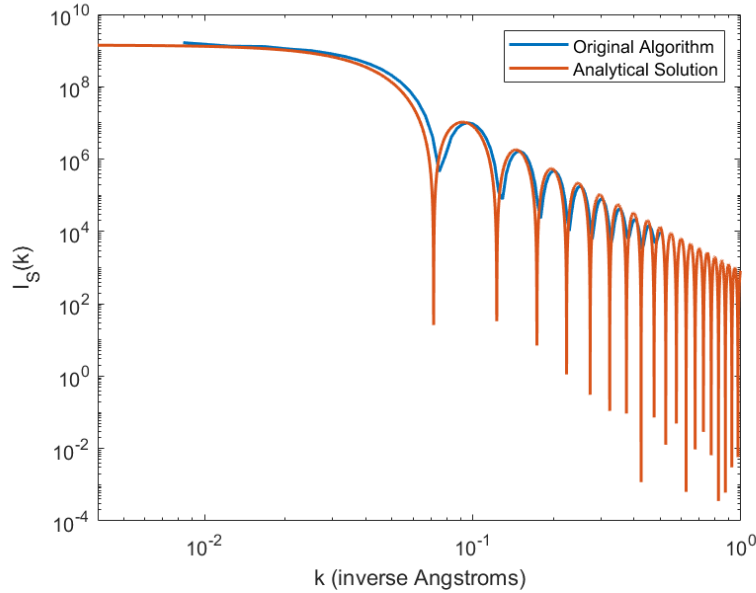


Figure 6: Original Orientation Averaging Algorithm vs Analytical Solution

Improved Algorithm

While the original algorithm offers a clear, step-by-step procedure for dis-

tributing scattering intensities across bins, its nested `for`-loop structure can be computationally expensive—especially for large grids where n_x or n_{xf} may exceed several hundred points. Furthermore, our analysis identifies two primary issues that lead to discrepancies between the algorithm’s output and the analytical predictions.

First Issue:

Rounding κ to the nearest integer bin index, followed by applying a linear “channel sharing” between adjacent bins, introduces subtle interpolation errors in the low- q region. For example, as shown in Fig. 6, the scattering intensity appears to be shifted to the right by approximately one bin. This shift is attributable to the design of the rounding function. To avoid division by zero, the algorithm initiates the sharing of scattering intensity from bin 1 rather than bin 0 (where $q = |\mathbf{q}| = 0$). Although this approach prevents division by zero, it results in a shifted scattering intensity once the correction $I'(q) = I(q)q^2$ is applied.

Second Issue:

The second problem also stems from the correction $I'(q) = I(q)q^2$. As depicted in Fig. 7, although the effective point counts as a function of q follow a quadratic trend for approximately 60% of the q range, this trend fails in the high- q region, leading to an underestimation of the scattering intensity in those bins. The effective counts are computed based on the weighted contributions of each point to each bin. Given the finite extent of the cube, the number of points within a spherical shell of radius q is not strictly proportional to q^2 when the radius approaches half the cube size. For q values beyond this threshold, fewer points are present since a portion of the spherical shell lies outside the cube.

To address these concerns, we propose an improved, more efficient method with three key modifications:

1. **Full-vectorization Approach:** Instead of iterating through all n_x^3 points using three nested loops, we first generate the entire 3D coordinate grid via `ndgrid` (or an equivalent function) to compute $|\mathbf{q}|$ values in a single vectorized operation. This mesh approach directly provides

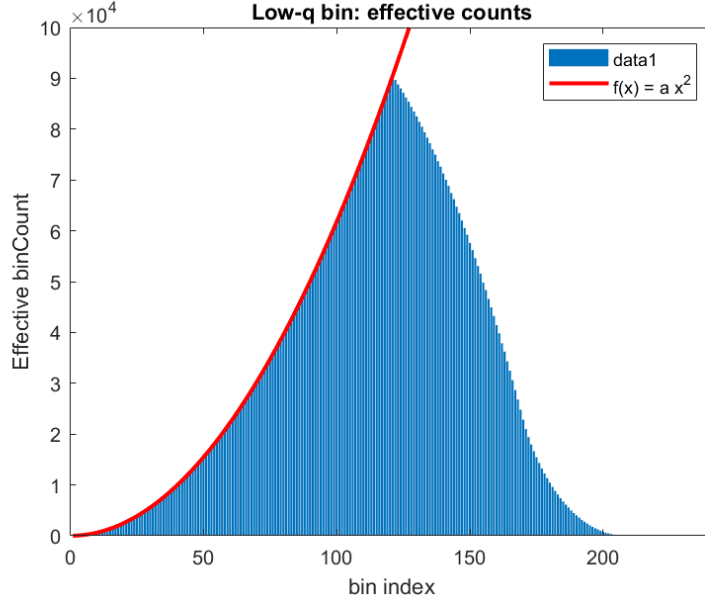


Figure 7: Effective Bin counts vs q . data1 denotes the effective counts of points falling into the bins along the q axis.

the spherical radius for each voxel in the frequency domain:

$$\mathbf{R} = \sqrt{X^2 + Y^2 + Z^2}$$

where (X, Y, Z) spans the 3D domain $\{-\frac{n_x}{2}, \dots, \frac{n_x}{2} - 1\}$ after shifting the zero-frequency component to the center. The corresponding intensities I_{q3D} and correction factors (sinc^2) are likewise extracted in vector form.

2. **Weighted Linear Interpolation via Accumulation:** Once each voxel's radius R_κ is known, we use a linear interpolation scheme to distribute its scattering intensity $I_{\text{DFT}}(\mathbf{q}_\kappa)$ across two adjacent bins. Denoting $\lfloor R_\kappa \rfloor$ by r_{floor} and the fractional part by $r_{\text{frac}} = R_\kappa - \lfloor R_\kappa \rfloor$, each voxel's intensity is apportioned to bin r_{floor} with weight $(1 - r_{\text{frac}})$ and to bin $(r_{\text{floor}} + 1)$ with weight r_{frac} . Mathematically:

$$I_{\text{scatt}}(r_{\text{floor}}) += I_{\text{DFT}}(\mathbf{q}_\kappa) (1 - r_{\text{frac}}), \quad I_{\text{scatt}}(r_{\text{floor}} + 1) += I_{\text{DFT}}(\mathbf{q}_\kappa) r_{\text{frac}}.$$

Crucially, we implement this “accumulation” in a vectorized manner

(e.g. via `accumarray` in MATLAB). This not only *eliminates* the triple nested loops but also ensures all voxel intensities are distributed in a single pass, dramatically reducing runtime.

3. **Better Normalization:** As in the original algorithm, we still apply the normalization by r^2 (or q^2) to account for the increasing number of points in concentric spherical shells at larger radii. However, as demonstrated above, this will cause the scattering intensities to be underestimated. To fix that, we normalize the scattering intensity with the effective counts of points instead of q^2 .

By combining these three modifications, the improved algorithm offers the following advantages:

- **Significantly Reduced Computational Cost:** The use of vectorized routines and direct accumulations avoids $\mathcal{O}(n_x^3)$ triple-loop overhead, thus scaling better to larger grids.
- **Better Physical Explanation:** The employment of a better normalization method avoids shifting the intensities intentionally, which makes more sense concerning physical meaning since there is scattering intensity at the center of the grid.
- **More Accurate Ripple Behavior:** Empirically, we observe better agreement with analytical scattering curves for spherical models, especially in the ripples, which was proved to be most significant for pattern recognition.

Figure 8 highlights these improvements. Compared to the original algorithm, the proposed scheme aligns more closely with the analytical result across all q ranges. As in the original approach, the high- q region still exhibits minor fluctuations due to finite sampling and boundary effects, but these are substantially reduced by using adaptive padding (Section 2.2.1) and continuous density correction (Section 2.2.2). Overall, this improved orientation averaging framework produces higher-quality SAXS curves in less time, thereby facilitating more efficient large-scale data generation for ML workflows [2, 3].

2.1.5 Data Pre-processing

After generating the simulated SAXS data for models of varying sizes and shapes, it is necessary to address the issue of inconsistent data lengths caused

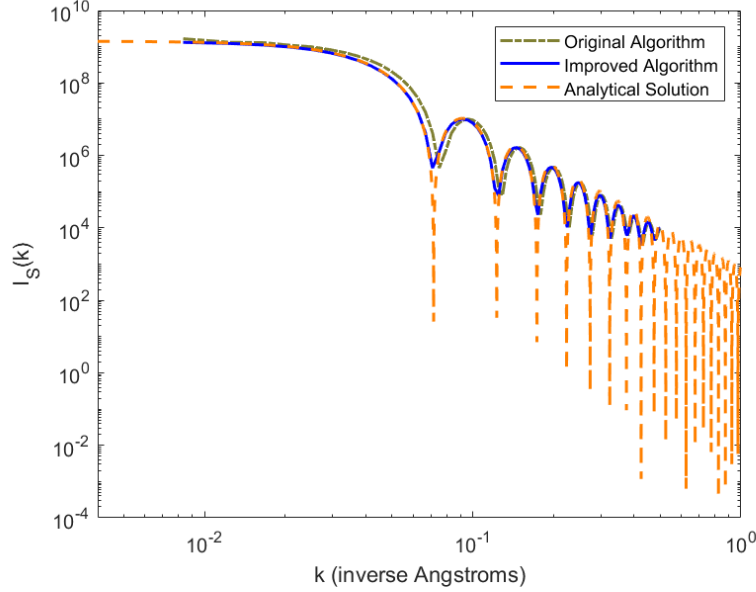


Figure 8: Comparison of the improved orientation averaging algorithm (red) with the original method (blue) and an analytical scattering curve (yellow) for a spherical LNP model.

by adaptive padding. Larger LNP models result in increased padding and a higher number of grid points (n_x), which leads to differences in the number of data points in the final SAXS curves. To overcome this, we adopt the solution proposed in [3], in which the scattering intensities are simulated for 500 (n_x) q values, uniformly distributed between:

$$q_{\min} \approx 0.04 \text{ nm}^{-1} \quad \text{and} \quad q_{\max} \approx 3.00 \text{ nm}^{-1},$$

These values were derived from an in-house experimental setup using an Anton Paar SAXSpoint 2.0 (Anton Paar, Graz, Austria) and represent a typical SAXS probing range while maintaining general applicability. Consequently, a global q range is defined for all LNP models. Although the LNPs are not assumed to be spherical, extrapolation using a Guinier function—defined as

$$I(q) = e^{A+Bq^2} \quad [4]$$

can provide a reasonable approximation for the scattering intensities of small LNPs before the Guinier region. To perform a weighted summation, we choose q_{\min} as the smallest q_{\min} value across the entire dataset. Mathemati-

cally, the universal q range is expressed as:

$$q_{\min} = \frac{2\pi}{n_x a} = 2\pi \times 0.00167 \text{ nm}^{-1}, \quad q_{\max} = \frac{\pi}{a} = \pi \times 0.5 \text{ nm}^{-1}.$$

With a universal q range established for all SAXS data, each dataset can be interpolated or extrapolated to conform to the same q axis with a fixed number of points (e.g., 500 points). Comprehensive experiments indicate that cubic spline interpolation yields superior results.

2.1.6 Data Synthesizing

This section outlines the methodology for dataset creation, as illustrated by the flow chart in Fig. 9. The process is divided into three sequential steps: Within-Type Data Augmentation, Cross-Type Data Synthesizing, and Poisson Noise Addition. At this stage, the ML algorithm is developed using two distinct types of LNPs (Type 3 and Type 4, as shown in Fig. 3). Due to computational constraints, 81 samples of 500-point SAXS curves have been generated, covering LNP sizes ranging from 20 nm to 100 nm in 1 nm increments. These samples form the basis for further augmentation, which is essential for mitigating fitting challenges and enhancing the model’s ability to capture relevant patterns in the data.

2.1.7 Within-Type Data Augmentation

In this study, the focus is on pure, dilute aqueous LNP solutions. Here, “pure” indicates that only LNPs are present in the sample, and “aqueous” signifies that water is the solvent. Given that water molecules contribute negligibly to the scattering intensity due to their small size relative to LNPs, we assume that only the LNPs contribute to the total scattering intensity. Moreover, as the solutions are dilute, interparticle interactions can be neglected. According to [4], the total scattering intensity of the sample can thus be treated as the weighted sum of the scattering intensities of individual particles, where each weight represents the number of particles of a specific size or type.

The first step in data augmentation is to increase the number of SAXS curves available for each LNP type. Based on the random residence time approach, the particle size distribution in finely divided systems is often observed to

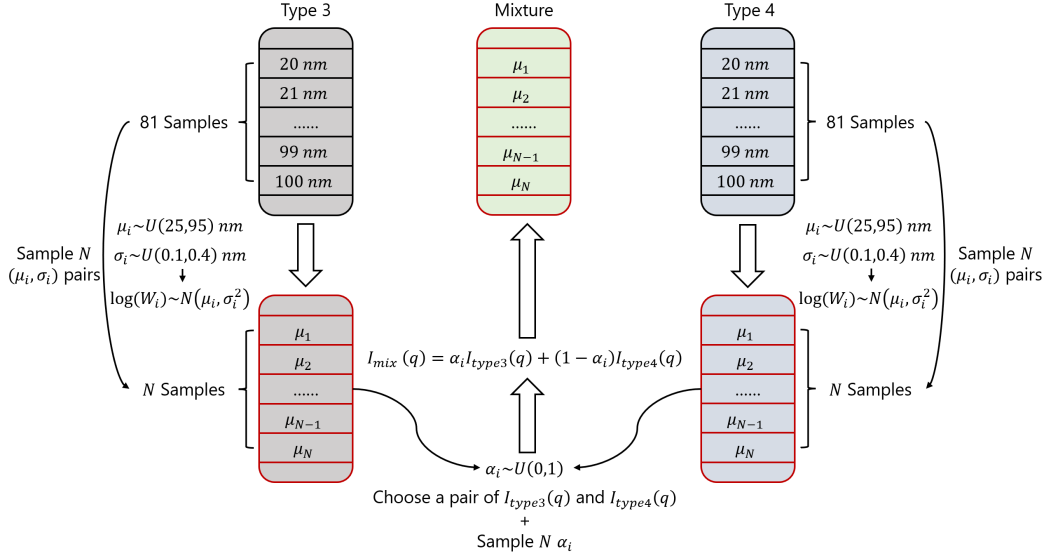


Figure 9: The Flow Chart of Data Synthesizing. Symbols: $U(a, b)$ denotes uniform distribution over interval (a, b) .

follow a lognormal distribution [5]. To capture this behavior, we uniformly sample a mean value, μ , within the range of 25 to 95, and a corresponding standard deviation, σ , within the range of 1 to 4. Each μ - σ pair defines a unique lognormal distribution. This distribution is divided into 81 intervals, from which 81 normalized weights (summing to 1) are extracted—each weight corresponding to one of the original samples. Fig. 10 illustrates the whole process by taking pure Type 3 LNP solutions as examples.

For each LNP type, the 81 original 500-point SAXS curves are multiplied by their respective weights, and the weighted intensities are summed to produce a new 500-point curve. By repeating this procedure N times, each LNP type is expanded into N samples. Each resulting SAXS curve is interpreted as the scattering result of a dilute LNP solution characterized by a size distribution determined by the sampled μ and σ . In this framework, μ serves as the label for the sample, while σ is regarded as a noise factor.

2.1.8 Cross-Type Data Synthesizing

To more accurately mimic real-world samples, heterogeneous LNP solutions are simulated via a weighted linear combination approach. Specifically, two

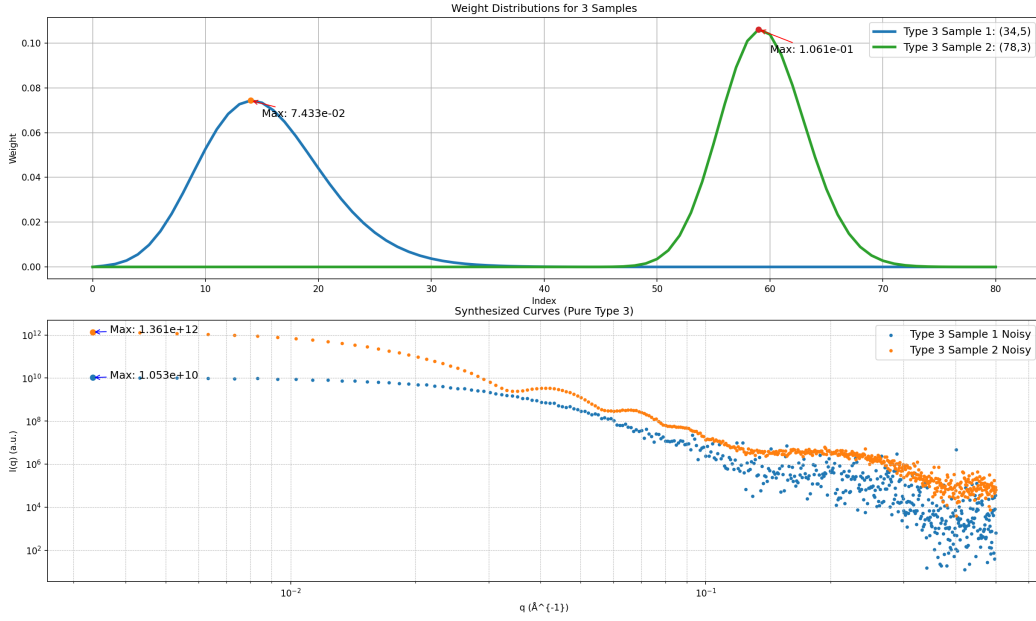


Figure 10: Demonstration of Within-Class Data Augmentation (With Poisson Noise Implemented)

samples are randomly selected from two different LNP types (e.g., Type 3 and Type 4), denoted as Curve 1 and Curve 2 with corresponding labels μ_1 and μ_2 . A random weight α is then drawn from a uniform distribution on the interval $(0, 1)$ to determine the relative contributions of each sample. The resulting heterogeneous SAXS curve is computed as:

$$\text{Curve}_{\text{heter}} = \alpha \times \text{Curve}_1 + (1 - \alpha) \times \text{Curve}_2,$$

with the associated label defined by the weighted average:

$$\mu_{\text{heter}} = \alpha \times \mu_1 + (1 - \alpha) \times \mu_2.$$

Here, α represents the fraction of one LNP type in the mixture, serving as a target for regression. In addition to predicting α , determining μ_{heter} provides further insights into the size distribution of the constituent LNP types. Both α and μ_{heter} are therefore recorded in the synthesized heterogeneous dataset.

2.1.9 Poisson Noise Addition

To accurately simulate the experimental measurement process, it is crucial to model the inherent noise originating from photon counting statistics [3, 6]. In SAXS experiments, the number of photons detected in each q -bin is governed by Poisson statistics, meaning that the variance of the measured intensity is proportional to its mean. However, when simulating theoretical SAXS curves, since it is impossible to know the devices' real-world calibration coefficients (if they exist), directly applying a Poisson noise model would produce unrealistic noise levels.

To address this, we adopt a lognormal noise model that ensures the added noise remains strictly positive and effectively handles multiplicative effects. For each theoretical SAXS curve $I(q)$, we introduce a multiplicative noise factor $N(q)$ defined as:

$$N(q) = \exp \left(\sqrt{\sigma_n^2(q)} \epsilon - \frac{\sigma_n^2(q)}{2} \right),$$

where $\epsilon \sim \mathcal{N}(0, 1)$ is a standard normally distributed random variable. This formulation leverages the properties of the lognormal distribution to guarantee that

$$E[N(q)] = 1,$$

thereby ensuring that the expected value of the noisy intensity remains $I(q)$, i.e.,

$$E[I_{\text{noisy}}(q)] = I(q).$$

The q -dependent variance $\sigma_n^2(q)$ is computed by:

$$\sigma_n^2(q) = \ln \left(1 + \frac{\alpha_n}{I(q)} \right),$$

where the noise scaling parameter α_n is sampled from a log-uniform distribution over a range (e.g., $[10^4, 10^{7.5}]$). This range is chosen based on extensive experimental evaluation to reflect the wide dynamic range observed in real SAXS measurements. Under this model, the variance of the noisy intensity approximates $\alpha I(q)$, consistent with the statistical behavior expected from

photon counting. Finally, the noisy SAXS curve is obtained by:

$$I_{\text{noisy}}(q) = I(q) \times N(q),$$

The examples of simulated SAXS curves are illustrated in Fig. 11 with a setting of $[10^5, 10^{8.5}]$ to provide more realistic results.

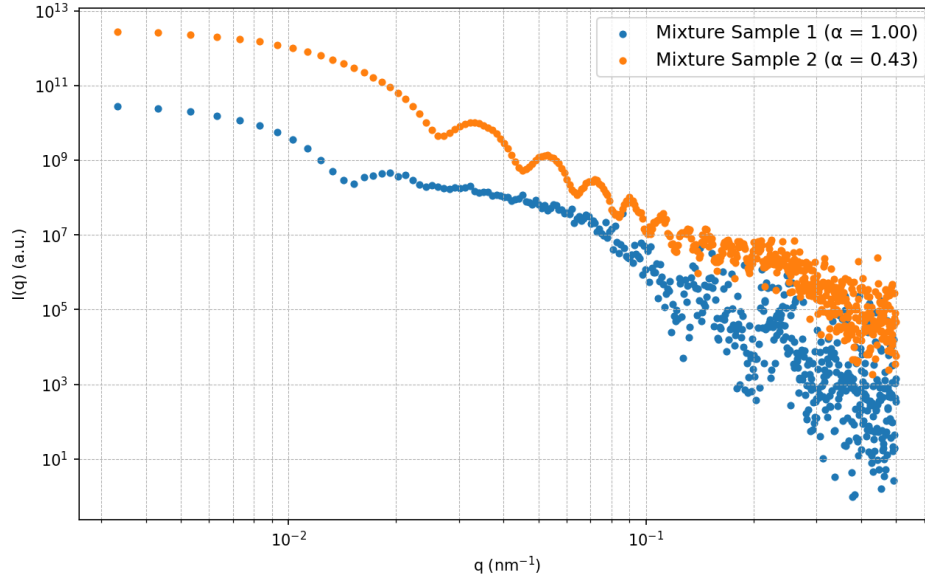


Figure 11: Example of simulated SAXS curves (With different mixing factor α)

2.2 ML Model Design

At this stage, two distinct convolutional neural network (CNN) architectures have been developed: one for classifying LNP solution types and another for regressing the fraction α . Both models incorporate effective normalization techniques—including batch normalization and feature scaling—to mitigate the effects of the large dynamic range present in the scattering intensity data, thereby enabling the networks to focus on the underlying shape of the SAXS curves.

The classification network (see Fig. 12) comprises two convolutional layers with 32 and 64 output channels and kernel sizes of 5 and 3, respectively. Each convolutional layer is followed by batch normalization and a ReLU activation

function. The feature maps are then flattened and passed through two fully connected layers, with a dropout layer (rate 0.3) inserted between them to reduce overfitting. A softmax function at the output layer generates a valid probability distribution over the target classes.

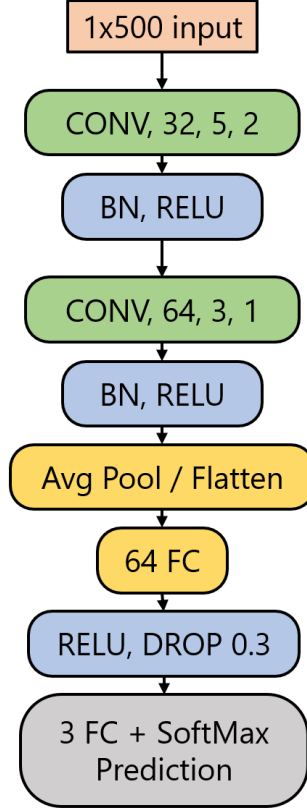


Figure 12: Network Architecture for Classification

The regression network (see Fig. 13) is designed to predict the mixing factor α and consists of three convolutional blocks. The convolutional layers in these blocks are configured with 32, 64, and 128 output channels and kernel sizes of 5, 5, and 3, respectively. Each layer is accompanied by batch normalization and a ReLU activation function. The final feature representation is then processed by three fully connected layers, with dropout layers (rate 0.15) inserted between them to mitigate overfitting. A sigmoid activation function at the final layer produces the regression output.

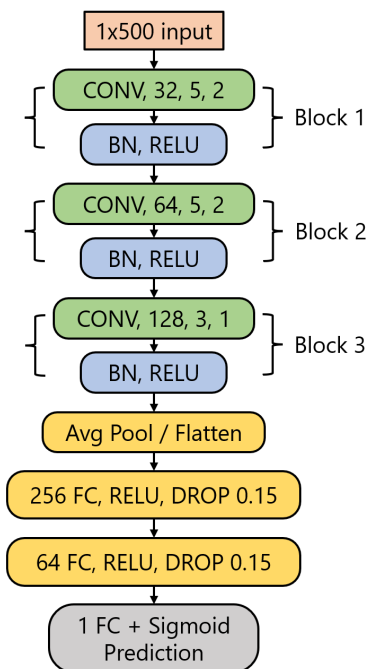


Figure 13: Network Architecture for Regression

3 Results

This section presents the experimental outcomes obtained from the developed data synthesizing pipeline and ML models. The results are organized into two subsections: one addressing the classification of LNP solution types and the other concerning the regression of the mixture fraction.

3.1 Results for Classification

The classification model demonstrated robust performance in differentiating between LNP solution types. As shown in Fig. 14, both the training and validation loss curves exhibit smooth convergence. The detailed classification metrics (see Table 1) indicate high precision and recall across all classes, with the F1-scores of 0.97 for Type 3, 0.98 for Type 4, and 0.93 for the Mixture category. An overall accuracy of 96% was achieved on a test set comprising 1200 samples. Furthermore, the confusion matrix in Fig. 15 reveals that misclassifications were minimal, thereby reinforcing the model’s discriminative capability in handling heterogeneous LNP solutions.

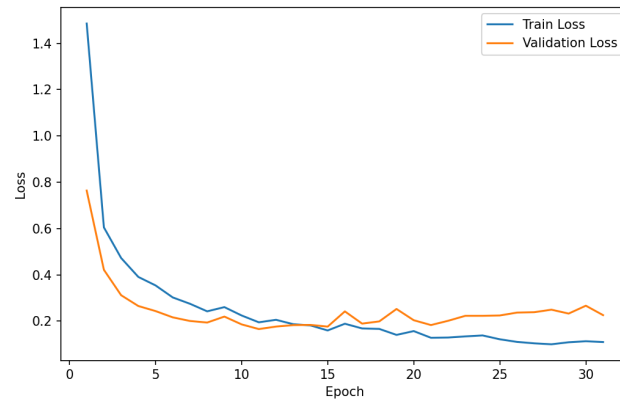


Figure 14: Training/Validation loss for classification

Table 1: Classification Report

Individual Classes				
Class	Precision	Recall	F1-score	Support
Type 3	0.95	0.99	0.97	437
Type 4	0.95	1.00	0.98	380
Mixture	0.98	0.89	0.93	383
Aggregate Metrics				
Accuracy	0.96			1200

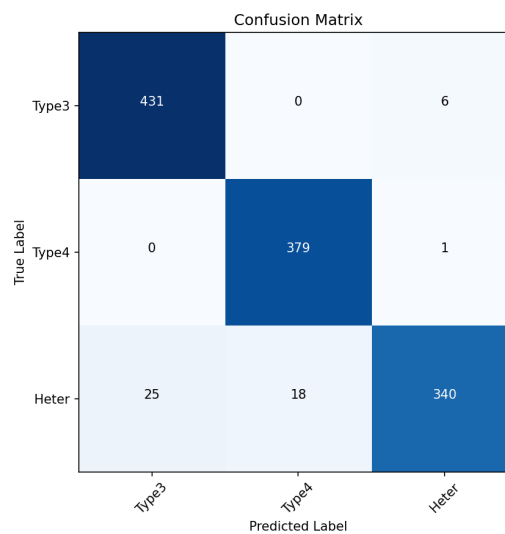


Figure 15: Confusion matrix for classification

3.2 Results for Fraction Regression

The regression model, designed to predict the fraction (α) in heterogeneous samples, exhibited excellent predictive performance. The model achieved an R^2 score of 0.9636 and a total test mean squared error (MSE) of 0.003019. The training losses are summarized in Fig. 16. These results underscore the model’s capacity to capture the underlying relationships in the SAXS data. Additionally, the regression performance is further corroborated by the log-scaled prediction results displayed in Fig. 17, which demonstrate a high degree of correlation between the predicted and true fraction values with minimal bias.

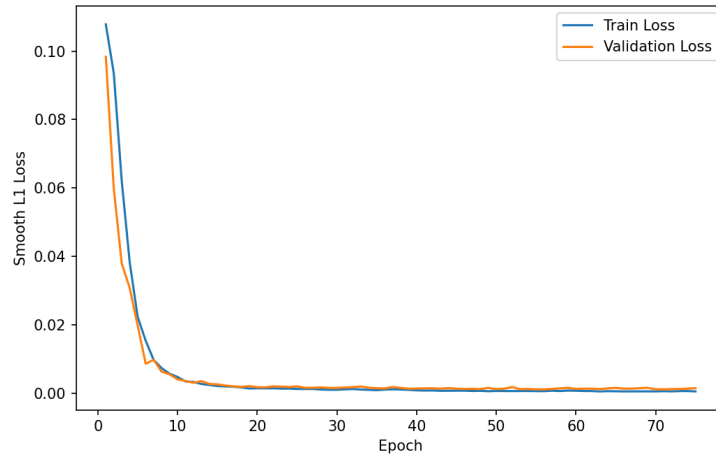


Figure 16: Training Losses for Regression Model

4 Conclusion

In this work, we presented a physics-based simulation pipeline and corresponding ML models for the characterization of LNP drug delivery vehicles using SAXS data. By carefully constructing realistic 3D LNP models, applying adaptive padding in the fast Fourier transform, introducing a continuous density correction factor, and incorporating lognormal and Poisson-like noise, we generated high-fidelity SAXS data across a broad range of particle sizes and mixture fractions.

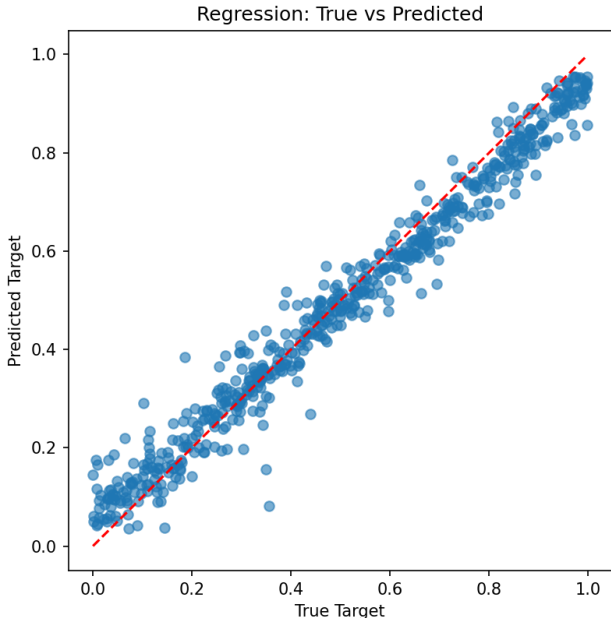


Figure 17: Prediction Results of Mixing Factor α

We then trained CNNs for two tasks: (1) classifying LNP solution types, and (2) predicting the mixture fraction in heterogeneous samples. The classification network achieved an overall accuracy of 96%, demonstrating its strong ability to distinguish between pure LNP solutions of different types and their mixtures. The regression model achieved a high R^2 of 0.9636, successfully capturing the underlying relationships in the data to predict the mixture fraction (α) of two LNP types.

These results highlight the potential of ML-driven SAXS analysis as a rapid and cost-effective complement or alternative to more resource-intensive methods such as cryo-EM. By offering quantitative insights into LNP size distributions and compositional fractions, the proposed approach paves the way for more efficient formulation development and quality control in nanomedicine. Nevertheless, future work is needed to further validate our pipeline with real experimental SAXS datasets, investigate the impact of sample polydispersity and other solvent effects, and extend the methodology to more complex multi-type LNP mixtures. Ultimately, coupling advanced simulation protocols with robust ML models holds promise for accelerating the understanding,

optimization, and deployment of LNP-based drug delivery systems.

References

- [1] A. Roesch, S. Zöls, D. Stadler, C. Helbig, K. Wuchner, G. Kersten, A. Hawe, W. Jiskoot, and T. Menzen, “Particles in biopharmaceutical formulations, part 2: An update on analytical techniques and applications for therapeutic proteins, viruses, vaccines and cells,” *Journal of Pharmaceutical Sciences*, vol. 111, pp. 933–950, 2022.
- [2] K. Schmidt-Rohr, “Simulation of small-angle scattering curves by numerical fourier transformation,” *Journal of Applied Crystallography*, vol. 40, pp. 16–25, 2007.
- [3] M. Röding, P. Tomaszewski, S. Yu, M. Borg, and J. Rönnols, “Machine learning-accelerated small-angle x-ray scattering analysis of disordered two- and three-phase materials,” *Frontiers in Materials*, vol. 9, p. 956839, 2022.
- [4] SASView, “Sasview documentation index: Corfunc technical documentation,” n.d., accessed: 2025-03-02. [Online]. Available: <https://www.sasview.org/docs/user/qtgui/Perspectives/Corfunc/corfunc-technical.html>
- [5] L. Kiss, J. Söderlund, G. Niklasson, and C. Granqvist, “The real origin of lognormal size distributions of nanoparticles in vapor growth processes,” *Nanostructured Materials*, vol. 12, no. 1–4, pp. 327–332, 1999.
- [6] S. Sedlak, L. Bruetzel, and J. Lipfert, “Quantitative evaluation of statistical errors in small-angle x-ray scattering measurements,” *Journal of Applied Crystallography*, vol. 50, no. 2, pp. 621–630, Mar 2017.