# SIMC2024 Challenge Report

Li Shiming        Liu Ye        Wang Boran

**Abstract**

This report details a computational approach to the Single Particle Imaging (SPI) challenge. We progress from simple geometric transformations for noiseless data to statistical clustering for unsupervised classification. Finally, to tackle the extreme noise and scale of Task 7, we introduce a **Stochastic Gradient Descent (SGD)** optimization method. By treating the unknown master image as a learnable weight matrix and employing a "Gradient Back-Rotation" technique, we successfully reconstruct the latent image using information from all 100,000 patterns simultaneously.

# Contents

# 1   Introduction

The core challenge of SPI is resolving the 3D structure of proteins from 2D diffraction patterns with unknown orientations. This challenge simulates this complexity through 2D pattern classification tasks ranging from noiseless, labeled data to massive, sparse, and extremely noisy datasets. Our methodology emphasizes data efficiency and scalable algorithms.

# 2   Phase I: Geometric Determination (Tasks 1 & 2)

## 2.1   Task 1 Approach (Noiseless Patterns)

Our approach relies on the discrete nature of the problem. Since there are only four possible orientations ($0°, 90°, 180°, 270°$) and the data is noiseless, we utilized a "Master Key" strategy. We designated the first image in the dataset as the reference ($Ref$). We computationally generated a dictionary of valid templates by rotating this reference:

$$Templates = \{\mathrm{Rot}(Ref, k \cdot 90°) \mid k \in \{0, 1, 2, 3\}\}$$

For every subsequent pattern $P$ in the dataset, we performed a pixel-wise equality check against these four templates. The match was exact and instantaneous, allowing us to categorize all 25 patterns with 100% accuracy without complex inference.

## 2.2   Task 2 Approach (Flattened Vectors)

Task 2 obscures the spatial relationships by flattening the $36 \times 36$ images into 1D vectors of length $L = 1296$. Since rotation is a 2D geometric operation, it cannot be directly applied to a 1D sequence. Our solution employs a topological transformation we call the "Fold and Unfold" trick:

1. **Fold:** We reshape the 1D vector back into a 2D matrix of shape $(36, 36)$.

2. **Rotate:** We apply the standard 2D rotation functions developed in Task 1.

3. **Unfold:** We flatten the rotated matrix back into a 1D vector.

By generating the 4 flattened templates of the reference vector, we reused the exact matching logic from Task 1. This proved that data representation (1D vs 2D) does not hinder classification as long as the transformation topology is known.

# 3   Phase II: Unsupervised Learning (Task 3)

## 3.1   Task 3 Approach (Unlabeled Scaling)

Scaling to 1,000 patterns with unknown labels (unsupervised learning) renders the "Master Key" strategy invalid, as we do not know which image is the correct $0°$ reference. We treated this as a clustering problem. We utilized **K-Means Clustering** with $K = 4$. We hypothesized that patterns with the same orientation would cluster together in the

high-dimensional pixel space ($D = 1089$). The algorithm minimizes the intra-cluster variance:

$$J = \sum_{j=1}^{4} \sum_{x \in S_j} ||x - \mu_j||^2$$

where $\mu_j$ is the centroid of cluster $j$. After convergence, the centroids ($\mu_j$) revealed the "average" structure of each orientation. By visually inspecting the sorted "Design Matrix" (Figure 5 in prompt), distinct bands of low entropy appeared, confirming that the algorithm successfully grouped the randomized data into four coherent orientation classes.

# 4    Phase III: Noise and Sparsity (Tasks 4, 5 & 6)

## 4.1    Task 4 Approach (Noisy Alignment)

In the presence of Poisson/Binomial noise, individual patterns are too sparse to visually identify. However, the signal persists statistically. We applied a "Consensus Alignment" strategy:

1. **Clustering:** We applied K-Means ($K = 4$) to group the noisy patterns. While individual samples are noisy, the cluster centroids (averages) suppress noise, revealing the underlying animal shape.

2. **Alignment:** The four centroids represent the animal at different orientations. To reconstruct the single master image, we treated Centroid 0 as the anchor. We computationally rotated Centroids 1, 2, and 3 to maximize their correlation with the anchor.

3. **Superposition:** Summing the aligned centroids further improved the Signal-to-Noise Ratio (SNR).

This confirmed that even when individual measurements are unreliable, the collective average of classified data yields a high-fidelity reconstruction.

## 4.2    Task 5: Likelihood Analysis

We derive the likelihood expressions and simplified ratios as follows:
(a) The likelihood for a rotation $r = 90°$ is defined as:

$$\begin{aligned}
\mathcal{L}(r = 90° \mid K, \mu) &\equiv \Pr(k_1 = 1 \mid \beta\lambda) \Pr(k_2 = 0 \mid \lambda) \Pr(k_3 = 0 \mid \beta\lambda) \Pr(k_4 = 0 \mid \beta\lambda) \\
&= \beta\lambda(1 - \lambda)(1 - \beta\lambda)^2
\end{aligned} \tag{1}$$

(b) Simplification of the ratio:

$$\text{Ratio} = \frac{\lambda(1 - \beta\lambda)^3}{\beta\lambda(1 - \lambda)(1 - \beta\lambda)^2} = \frac{1 - \beta\lambda}{\beta(1 - \lambda)} \tag{2}$$

(c) Derivation of $\beta$:

$$\text{Since} \quad \frac{1 - \beta\lambda}{\beta - \beta\lambda} = 1 \implies 1 - \beta\lambda = \beta - \beta\lambda \implies \beta = 1 \tag{3}$$

(d) & (e) Further simplifications for higher powers:

$$\text{(d)} \quad \frac{\lambda(1-\beta\lambda)^{15}}{\beta\lambda(1-\lambda)(1-\beta\lambda)^{14}} = \frac{1-\beta\lambda}{\beta-\beta\lambda} \tag{4}$$

$$\text{(e)} \quad \frac{\lambda^3(1-\beta\lambda)^{13}}{(\beta\lambda)^3(1-\lambda)^3(1-\beta\lambda)^{10}} = \left(\frac{1-\beta\lambda}{\beta-\beta\lambda}\right)^2 \tag{5}$$

(f) **Log-Likelihood derivation:**
For the first $M$ pixels where $\mu = \lambda$:

$$\ln(\mathcal{L}(\text{aligned} \mid \vec{\mu}, \vec{k})) = \sum_{i=1}^{M}\left(k_i \ln(\lambda) + (1-k_i)\ln(1-\lambda)\right) \tag{6}$$

Since $k_i$ follows a binomial distribution, we substitute the expected value $\lambda$:

$$\text{Term}_1 = M\left(\lambda\ln(\lambda) + (1-\lambda)\ln(1-\lambda)\right) \tag{7}$$

For the last $N - M$ pixels where $\mu = \beta\lambda$:

$$\text{Term}_2 = (N-M)\left(\beta\lambda\ln(\beta\lambda) + (1-\beta\lambda)\ln(1-\beta\lambda)\right) \tag{8}$$

Thus, the total average log-likelihood is:

$$\begin{aligned}
\text{Avg. Log-Likelihood} = {} & M\lambda\ln(\lambda) + M(1-\lambda)\ln(1-\lambda) \\
& + (N-M)\beta\lambda\ln(\beta\lambda) + (N-M)(1-\beta\lambda)\ln(1-\beta\lambda)
\end{aligned} \tag{9}$$

## 4.3   Task 6 Approach (Sparse Engineering)

Task 6 involves scaling to $6.5 \times 10^4$ patterns, making memory management critical. The data is "sparse," meaning most pixels are zero. Our approach focused on engineering optimization:

1. **Sparse Storage:** We utilized **the Compressed Sparse Row (CSR)** format, storing only non-zero indices. This reduced memory footprint by over 90%.

2. **Mini-Batch Clustering:** Standard K-Means is slow on large datasets. We used '**MiniBatchKMeans**', which updates centroids using small, random batches of data rather than the entire dataset.

This allowed us to process the massive dataset on standard hardware while maintaining the classification accuracy achieved in Task 4.

# 5   Phase IV: The Neural Solution

## 5.1   Task 7: Two Solvers for the Inverse Problem

Task 7 presents the ultimate challenge: $10^5$ patterns with extreme noise. The hint suggests we must "share data amongst orientations." We implemented and compared two distinct algorithmic strategies to solve this.

### 5.1.1 Approach A: Iterative Alignment (Expectation-Maximization)

Our first approach is a heuristic "Alternating Minimization" or EM-style algorithm. It solves the "Chicken and Egg" problem of SPI: we need the master image to determine orientations, but we need orientations to reconstruct the master image.

1. **Initialization:** We start with a "Guess Model" (usually the global average).

2. **Alignment (E-Step):** We rotate the Guess Model to 4 positions and match every pattern in the dataset to the closest rotation.

3. **Reconstruction (M-Step):** Once patterns are assigned to an orientation, we back-rotate them to the 0° frame and average them to create a refined Guess Model.

4. **Iteration:** We repeat this cycle. As the model improves, the alignment becomes more accurate, which in turn further improves the model.

### 5.1.2 Approach B: Stochastic Gradient Descent (SGD)

Our second approach reframes the task as a **Supervised Optimization Problem**. Instead of averaging, we define the Master Image as a learnable weight matrix $\mathbf{W}$ and minimize the reconstruction error using SGD.

- We calculate the gradient of the error for each batch.

- Crucially, we employ a **"Gradient Back-Rotation"** technique: if a pattern matches the master rotated by 90°, its gradient is back-rotated by −90° before updating $\mathbf{W}$.

- This allows information to flow from all samples to the canonical model simultaneously.

## 5.2 Iterative Reconstruction

We present two algorithms. Algorithm 3 is stable and intuitive, while Algorithm 4 is scalable and optimization-based.

---

**Algorithm 1** Iterative Alignment (EM-Style)

---

1: **Input:** Dataset $\mathbf{X}$
2: **Initialize:** Model $\mathbf{M} \leftarrow \mathrm{Mean}(\mathbf{X})$
3: **for** $iteration = 1$ to 5 **do**
4:     $\mathbf{M}_{rots} \leftarrow \mathrm{GenerateRotations}(\mathbf{M})$
5:     $NewModel \leftarrow \mathbf{0}$
6:     **for** each pattern $p$ in $\mathbf{X}$ **do**
7:                                                                     ▷ Find best alignment
8:         $k^* \leftarrow \mathrm{argmin}_k ||p - \mathbf{M}_{rots}[k]||$
9:                                                          ▷ Accumulate back-rotated pattern
10:         $NewModel \leftarrow NewModel + \mathrm{Rot}(p, -k^*)$
11:     **end for**
12:     $\mathbf{M} \leftarrow NewModel/N_{samples}$
13: **end for**

---

**Algorithm 2** SGD with Gradient Back-Rotation
___
1: **Input:** Dataset $\mathbf{X}$, Learning Rate $\eta$
2: **Initialize:** $\mathbf{M} \leftarrow \text{Mean}(\mathbf{X})$
3: **for** $epoch = 1$ to $N$ **do**
4:     Shuffle $\mathbf{X}$
5:     **for** each batch $\mathbf{B}$ in $\mathbf{X}$ **do**
6:        ... (Same as previous SGD logic) ...
7:        $Grad_{aligned} \leftarrow \text{Rot}(Grad, -k^*)$
8:        $\mathbf{M} \leftarrow \mathbf{M} + \eta \cdot \text{Mean}(Grad_{aligned})$
9:     **end for**
10: **end for**
___

## 5.3 Reconstruction via Iterative Alignment & SGD

Task 7 demonstrated the power of our algorithmic evolution.

- **Iterative Alignment (EM):** This method proved highly stable. Within just 3 iterations, the shapeless global mean converged into a recognizable face. It validated our hypothesis that "alignment improves reconstruction, and reconstruction improves alignment."

- **SGD Optimization:** The neural-style approach achieved similar high-fidelity results but offered greater flexibility in batch processing. The "Gradient Back-Rotation" mechanism successfully accumulated signal while cancelling out noise.

Although without datas ,we can really kown which one is better and could get a more clear picture,by analyzing methods,we can find that:
Both methods converged to the same solution: a clear portrait of a person (likely a famous scientist), validating that our approaches are mathematically robust.

# 6 Conclusion

We approached these tasks with a philosophy of simplicity evolving into optimization.

1. **Geometry:** Tasks 1-2 solved by geometric transformation.

2. **Statistics:** Tasks 3-6 solved by unsupervised clustering.

3. **Algorithm Design:** For Task 7, we demonstrated that the SPI "inverse problem" can be solved via two powerful paths: **Iterative Expectation-Maximization** and **Stochastic Gradient Descent**. Both methods effectively "share data" across orientations to recover the latent signal.

# 7 AI Use Report

**Trae IDE** was used for inline code completions and debugging.
**Gemini 3.0** was used to generate LaTeX code and improve the SGD code for Task 7.
WARNING:
All the codes and the paper were done without data,so the paper is lack of pictures and

all the data analysis were done **only** through Maths but without experimenting. The codes assume that there is a file called "endeavor.npz".