# Identifying Structural Properties of Proteins from X-ray Free Electron Laser Diffraction Patterns

**Paula Olaya***, Silvina Caino-Lores*, Vanessa Lama*, Ria Patel*, Ariel Keller Rorabaugh*, Osamu Miyashita[ǂ], Florence Tama[ǂ†], and Michela Taufer*
*University of Tennessee, Knoxville, USA
[ǂ]Center for Computer Science, RIKEN, Kobe, Japan
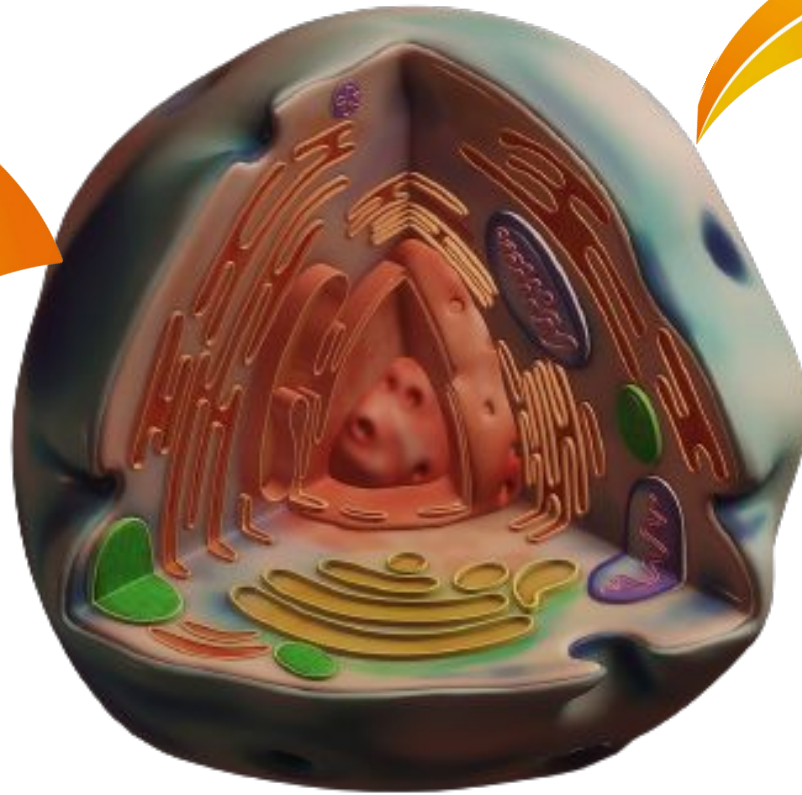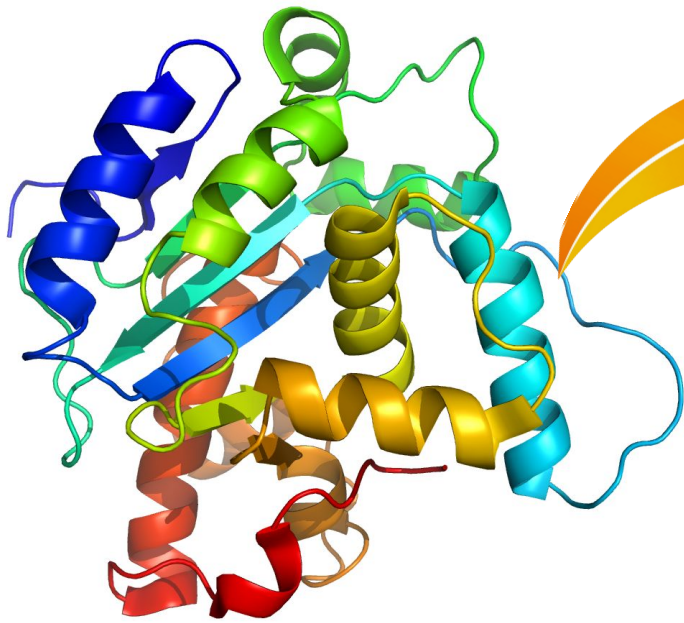[†]Nagoya University, Nagoya, Japan
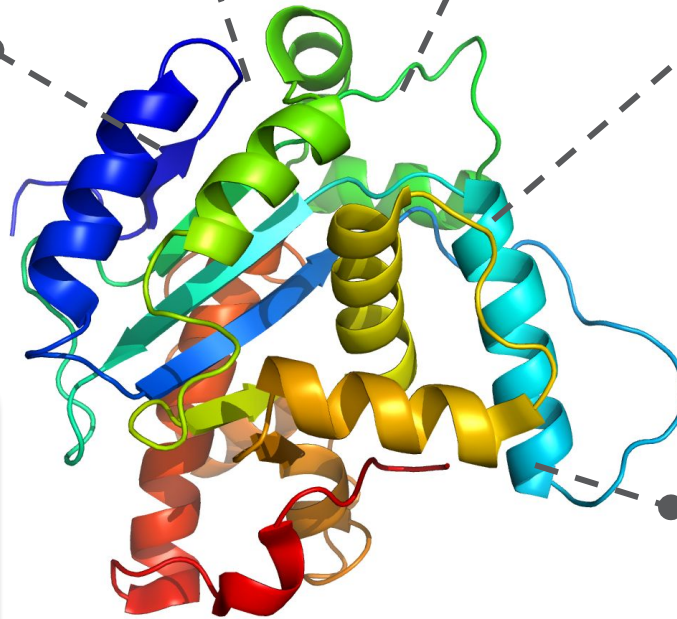
# Proteins are essential for all living organisms

**Catalysis:** Increasing the rate of a chemical reaction within cells

**Transportation:** Moving materials within a cell and the organism

**Structure:** Providing structure and support for cells

**Signaling:** Receiving, processing, and transmitting signals within the cell and with the environment
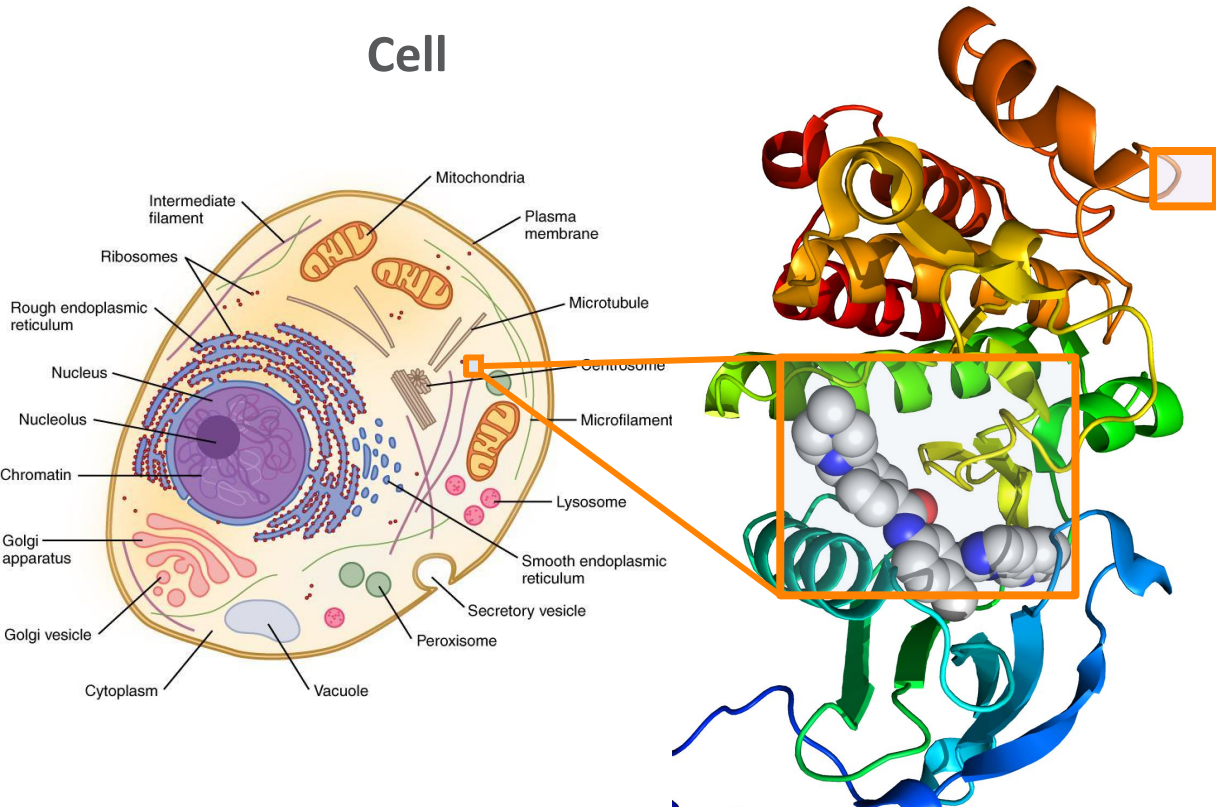
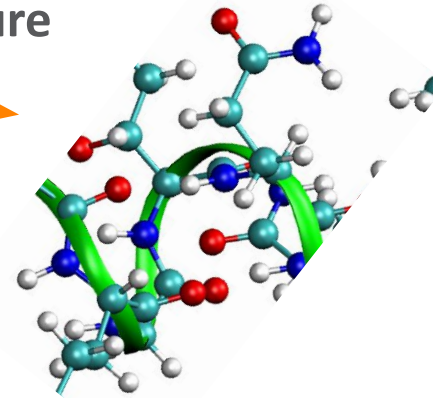**Proteins are responsible for many vital cellular functions**

**Antibodies:** Helping to protect the body from foreign particles, such as viruses and bacteria

- Structural biology explains 3D structures of biomolecules
- Biomolecules are **proteins**/RNA/DNA
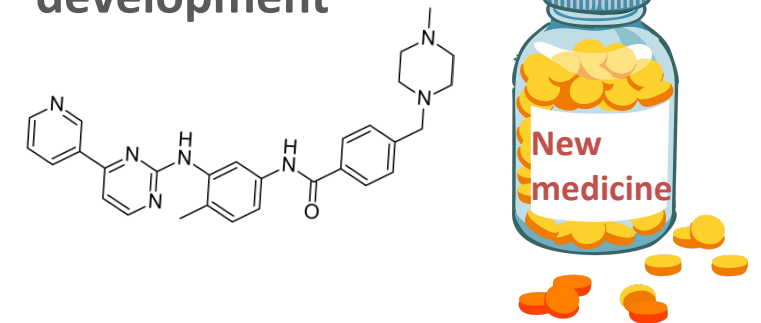
**Cell**



**Protein 3D structure**

- Information on atomic positions
- Hundreds of thousands of atoms and more

Gleevec®
Anticancer drug
(Leukemia)

**Drug development**

New medicine

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# There are experimental methods (X-ray, cryo-EM, SAXS, XFEL) to obtain partial information about the 3D protein structure
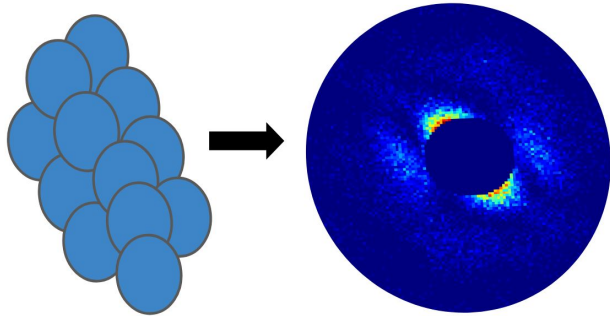
X-ray Free Electron Laser (XFEL) beams create 2D diffraction patterns that reveal properties of the 3D protein structure
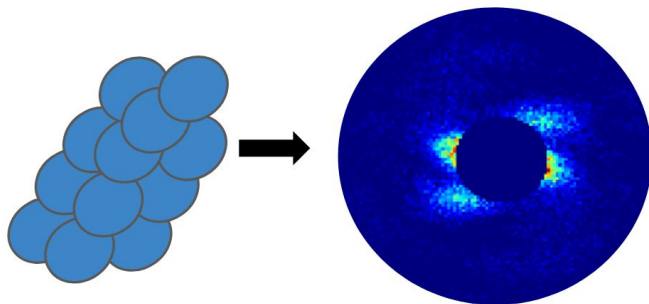
# Structural properties: Orientation



**Orientation 1**

Φ, θ, Ψ = 24°, 151°, 346°



**Orientation 2**

Φ, θ, Ψ = 145°, 128°, 291°

Orientation refers to the **placement of the incident beam** with respect **to a protein structure**

- Φ (Azimuth) = [-180,180]
- Θ (Altitude) = [0,180]
- Ψ (Psi or rotation angle) = [0, 360]

THE UNIVERSITY OF TENNESSEE KNOXVILLE

# Structural properties: Conformation

Conformation is the **shape adopted by a protein** and is caused by the rotation of the protein atoms around one or more single bonds

**Conformation A**

$\Phi, \theta, \Psi = 24^\circ, 151^\circ, 346^\circ$



**Conformation B**

$\Phi, \theta, \Psi = 34^\circ, 139^\circ, 106^\circ$

# Structural properties: Protein type

Protein type refers to the type and number of amino acids composing a protein

- 20 different type of amino acids

- Amino acids can combined in different ways to make a protein
  - sequence
  - number (up to thousands)

**Protein type A**
**Conformation A**
$\Phi, \theta, \Psi = 24^o, 151^o, 346^o$

**Protein type B**
**Conformation C**
$\Phi, \theta, \Psi = 84^o, 32^o, 82^o$

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# Structural properties

**Orientation 1**
Φ, θ, Ψ = 24º, 151º, 346º



**Conformation A**
Φ, θ, Ψ = 24º, 151º, 346º



**Protein type A**
**Conformation A**
Φ, θ, Ψ = 24º, 151º, 346º



**Orientation 2**
Φ, θ, Ψ = 145º, 128º, 291º



**Conformation B**
Φ, θ, Ψ = 34º, 139º, 106º



**Protein type B**
**Conformation C**
Φ, θ, Ψ = 84º, 32º, 82º

Experimental 2D diffraction patterns

3D Structure (Fourier Space)

3D Structure (Real Space)

XFEL

Identifying the structural properties embedded in the **2D diffraction pattern** is key for the 3D reconstruction and understanding the protein's structure

THE UNIVERSITY OF TENNESSEE KNOXVILLE

We need to **integrate** the **experimental** methods with **computational** frameworks **to gain information on structure and dynamics** and accelerate scientific discovery



$a_1 = -10.1$
$a_2 = 151.3$
$a_3 = 305.8$
$conf = 1n0u$
$pt = EF2$

Our goal is to design and implement a **ML-based framework** that predicts simultaneously the three structural properties from protein diffraction patterns

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# Framework design consideration 1

1. Simultaneous multi-output and multi-type predictions

**Orientation**

| **Continuous values** |
| --- |
| Angle 1 = [-180,180] |
| Angle 2 = [0,180] |
| Angle 3 = [0,360] |

**Conformation**

| **Categorical values** |
| --- |
| Conf. A1 or Conf. A2 or … or Conformation NN |

**Protein type**

| **Categorical values** |
| --- |
| Protein A or Protein B or … or Protein N |

# Framework design consideration 1

1. Simultaneous multi-output and multi-type predictions



$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{bmatrix}$

**Feature vectors**

**ML Model**

**Orientation**

**Continuous values**
Angle 1 = [-180,180]
Angle 2 = [0,180]
Angle 3 = [0,360]

**Conformation**

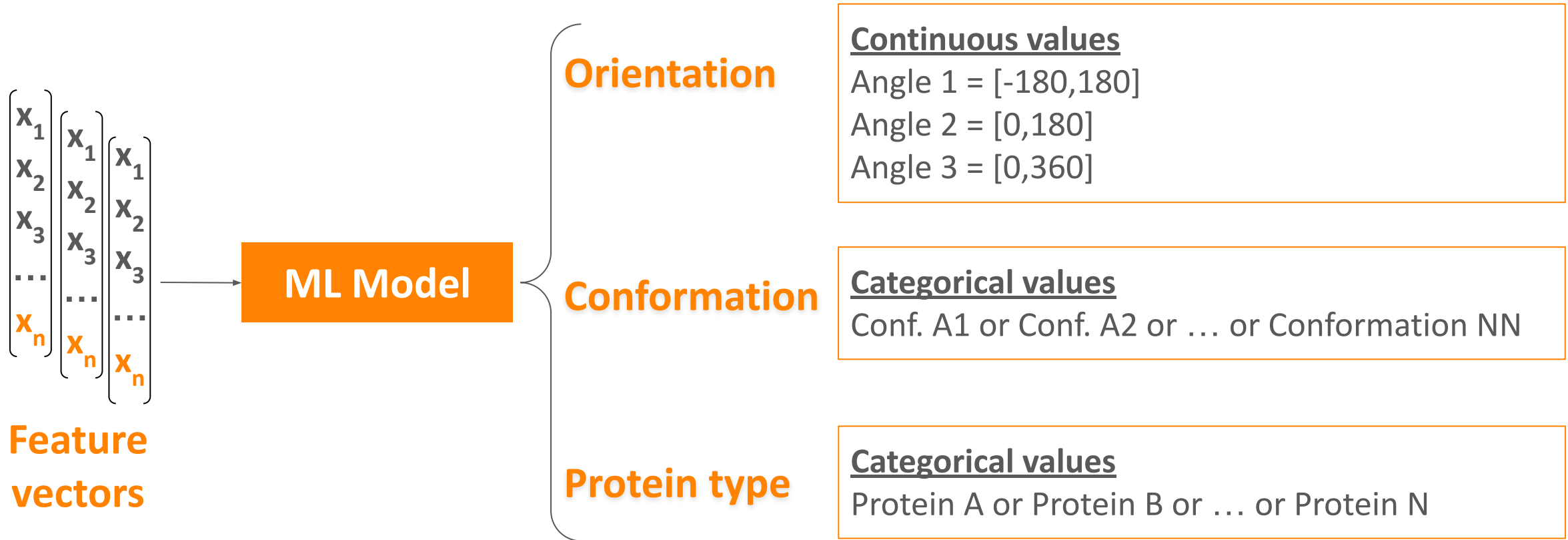**Categorical values**
Conf. A1 or Conf. A2 or … or Conformation NN

**Protein type**

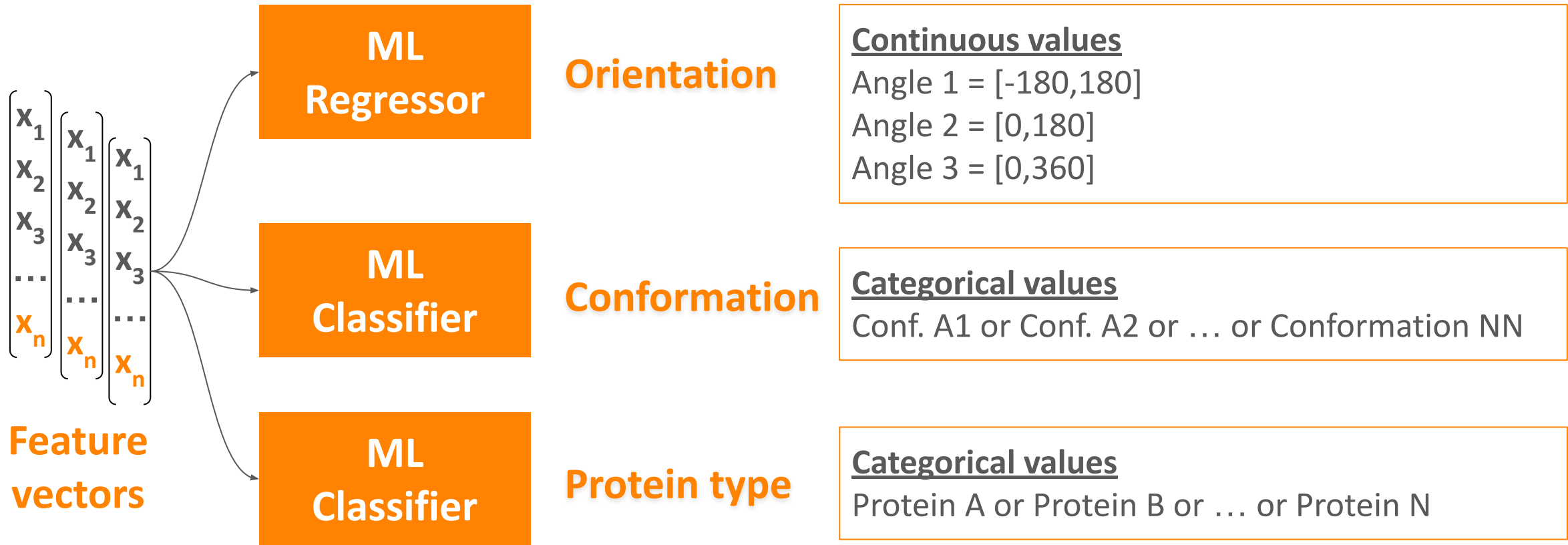**Categorical values**
Protein A or Protein B or … or Protein N

THE UNIVERSITY OF TENNESSEE KNOXVILLE

# Multiple ML models

We define three different ML models for the three predictions



$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{pmatrix}$$

**Feature vectors**

**ML Regressor** — **Orientation**

**Continuous values**
Angle 1 = [-180,180]
Angle 2 = [0,180]
Angle 3 = [0,360]

**ML Classifier** — **Conformation**

**Categorical values**
Conf. A1 or Conf. A2 or … or Conformation NN

**ML Classifier** — **Protein type**

**Categorical values**
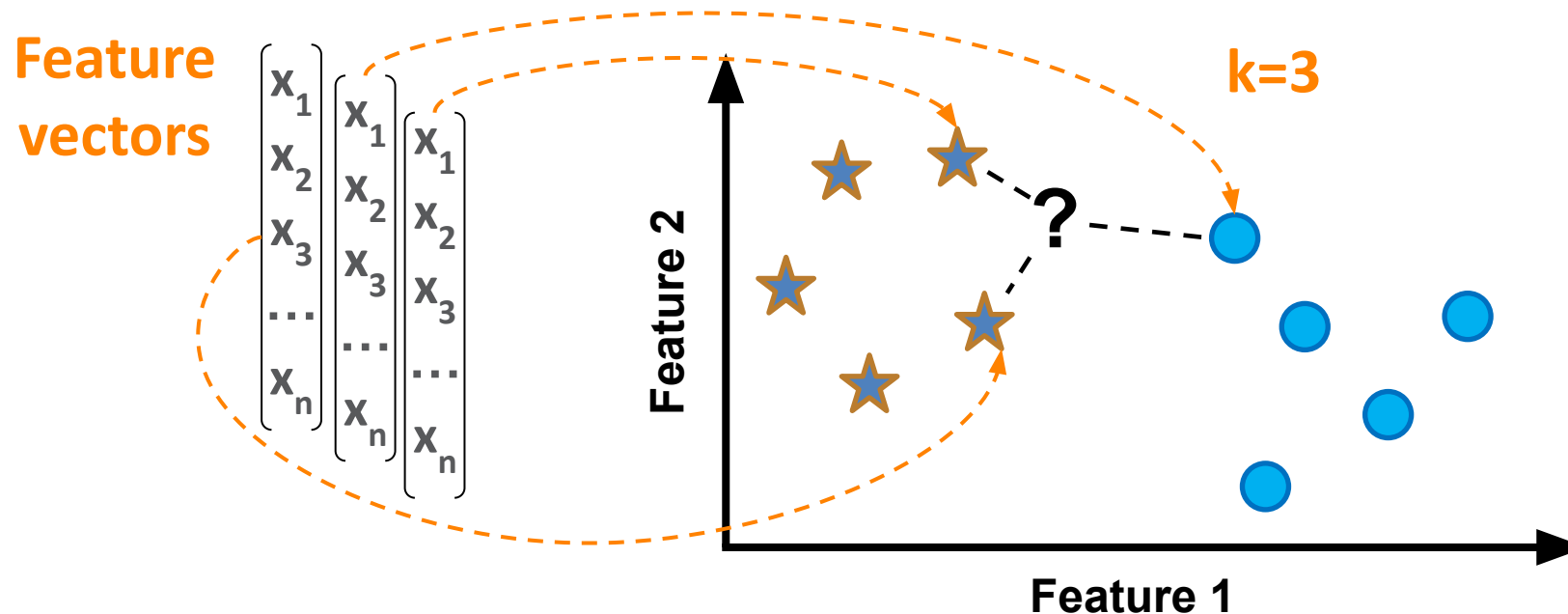Protein A or Protein B or … or Protein N

THE UNIVERSITY OF TENNESSEE KNOXVILLE

# Multiple kNN models

We select kNN (k-Nearest Neighbors) because of it high accuracy in both classification and regression problems and low execution costs

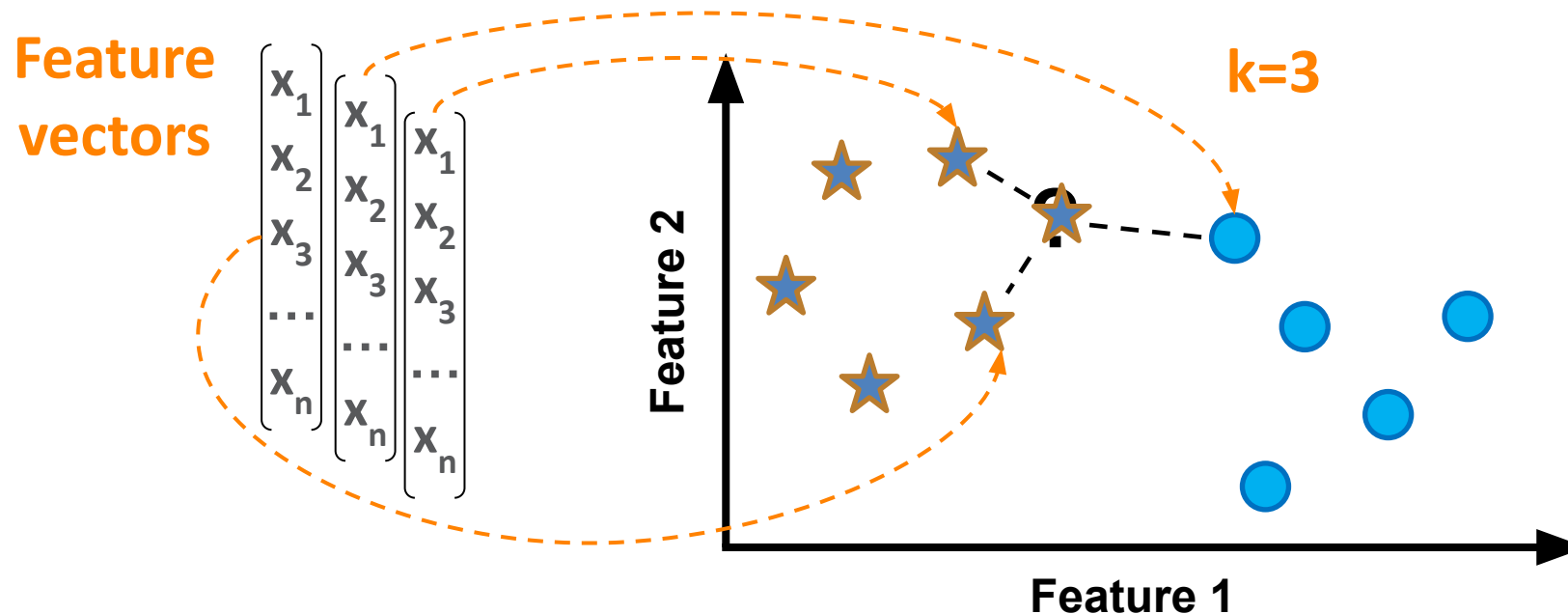THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# k-Nearest Neighbors (kNN)

This algorithm looks at the *K* nearest neighbors of a new data point (in feature space) to determine the predicted value
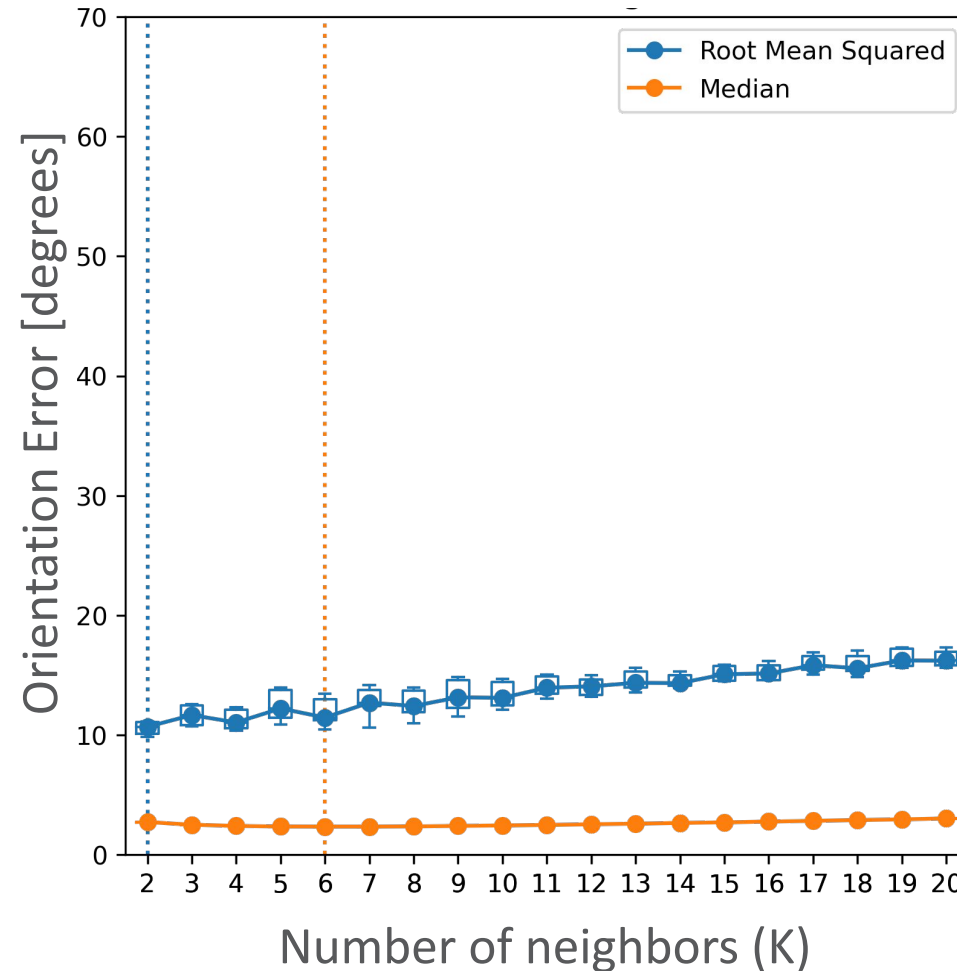
# k-Nearest Neighbors (kNN)

This algorithm looks at the *K* nearest neighbors of a new data point (in feature space) to determine the predicted value
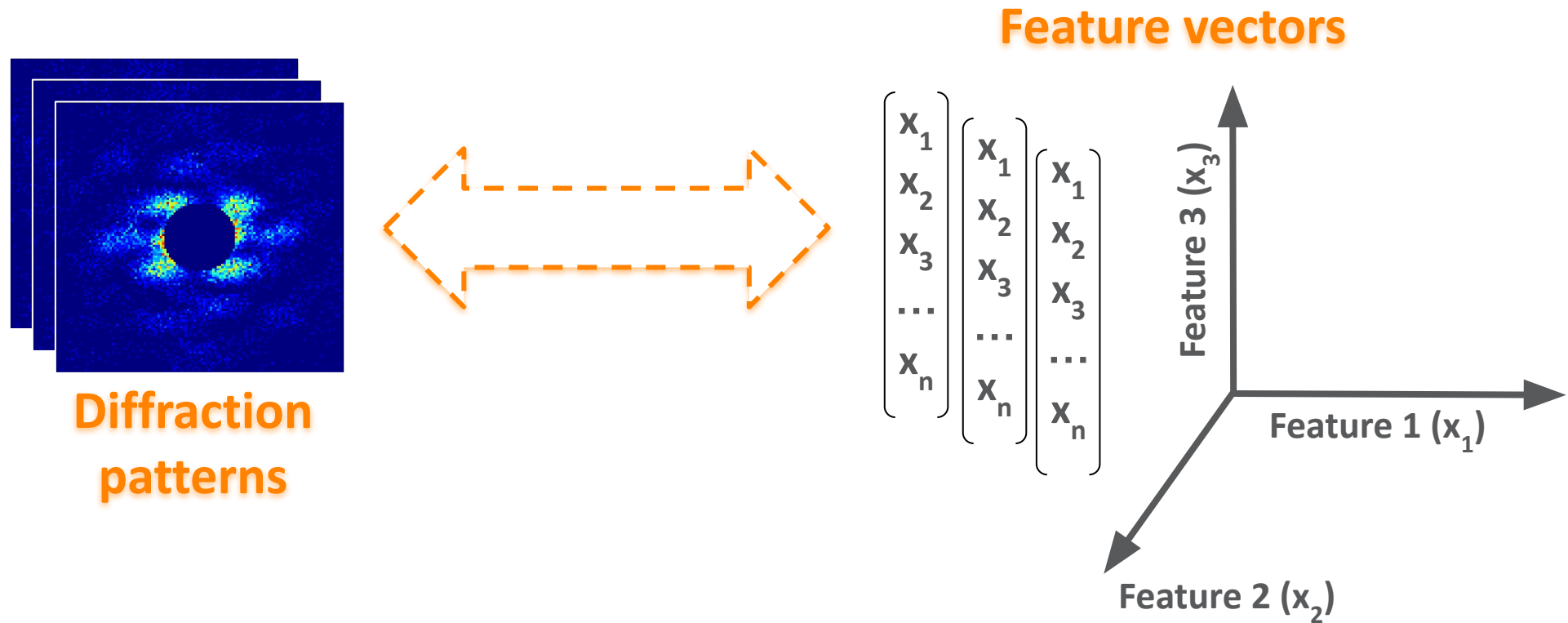
# Selecting the K number of neighbors

The K number of neighbors is critical for the prediction



An analysis of the root mean square error (RMSE) of the degree allows our framework to identify the most suitable K number of neighbors
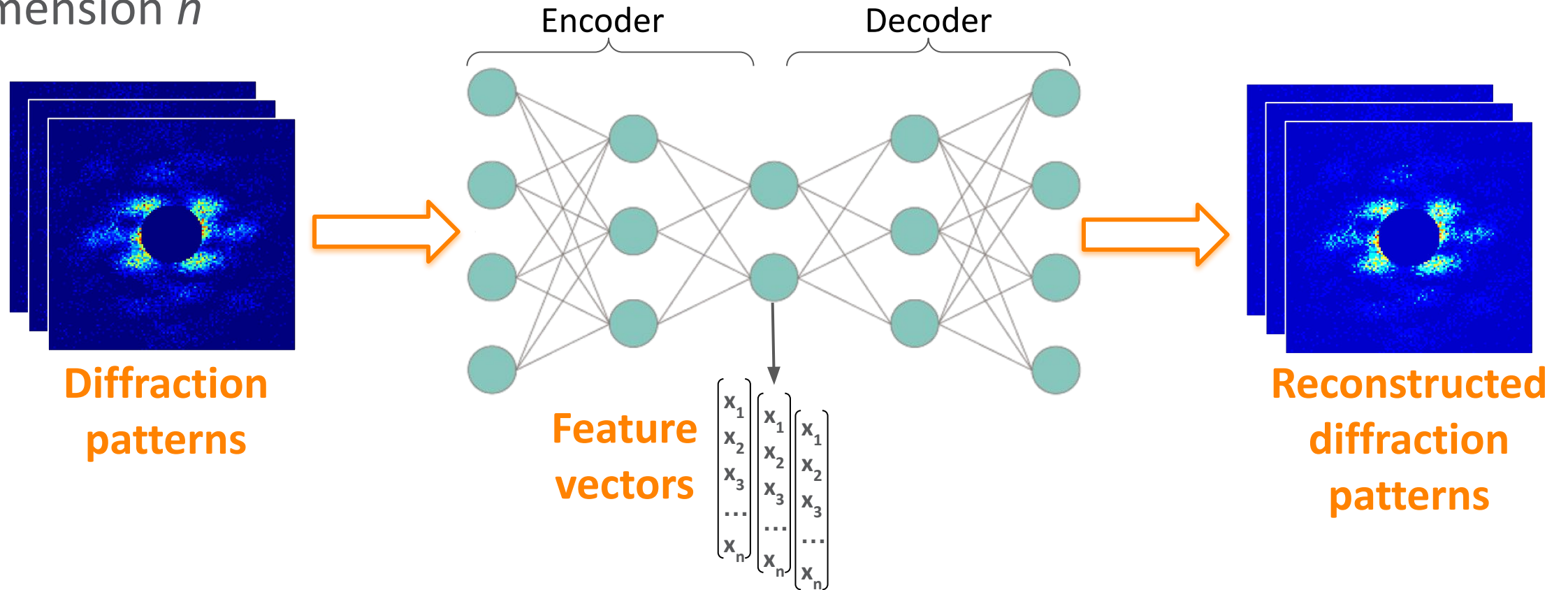
THE UNIVERSITY OF TENNESSEE KNOXVILLE

# Framework design consideration 2

2. Transformation from diffraction patterns to feature vectors

# From diffraction patterns to feature vectors

We use an **autoencoder** to represent diffraction patterns in feature vectors of dimension *n*
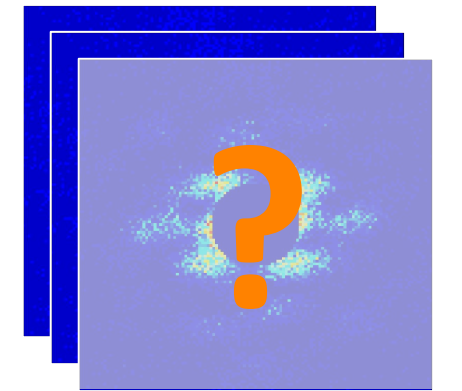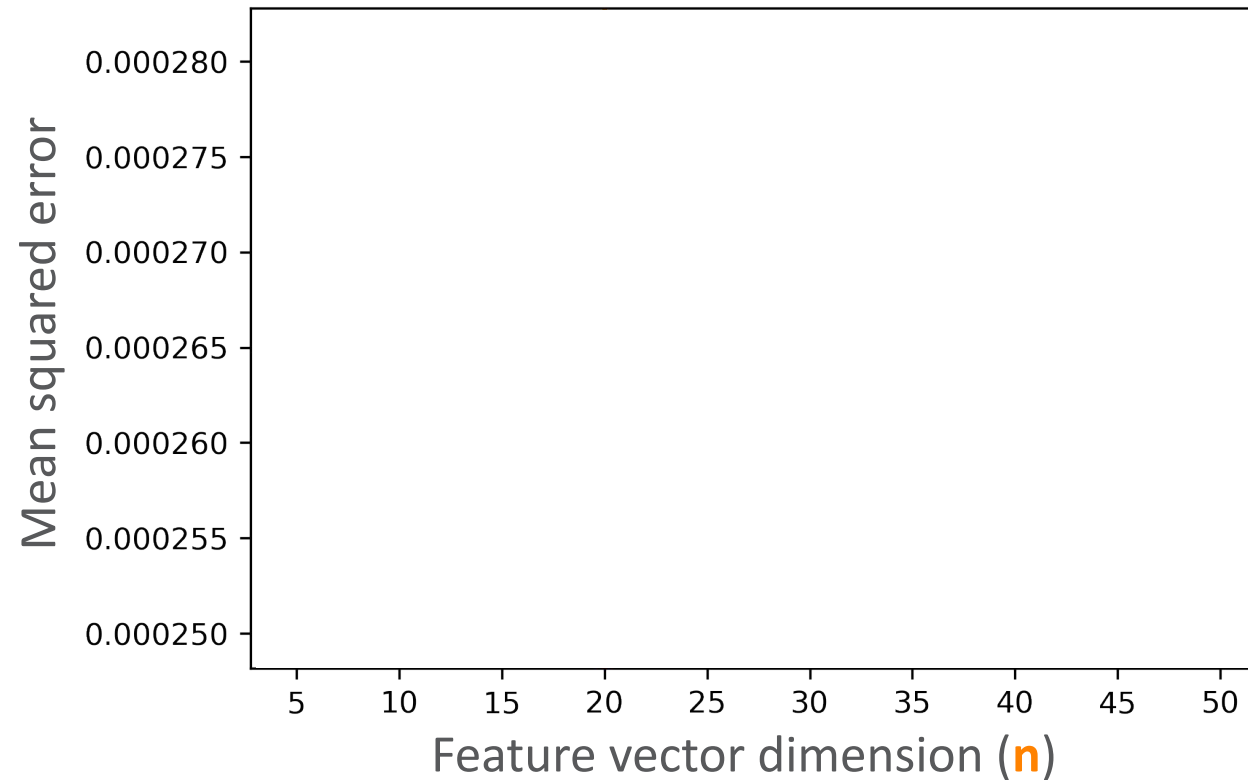


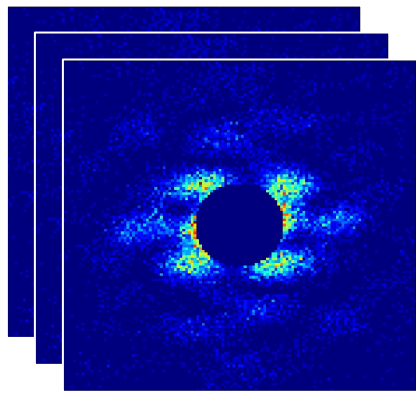Encoder

Decoder

**Diffraction patterns**

**Feature vectors** $\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{pmatrix}$

**Reconstructed diffraction patterns**

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# Feature vectors dimension (n)

The dimension of the feature vector has to be sufficient to faithfully reconstruct the original diffraction patterns
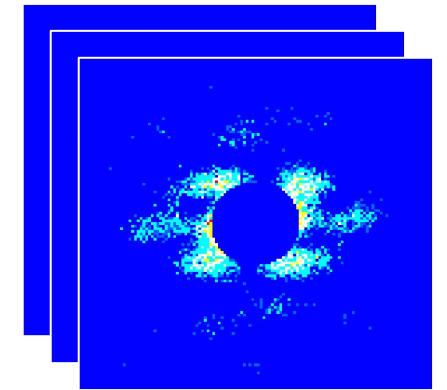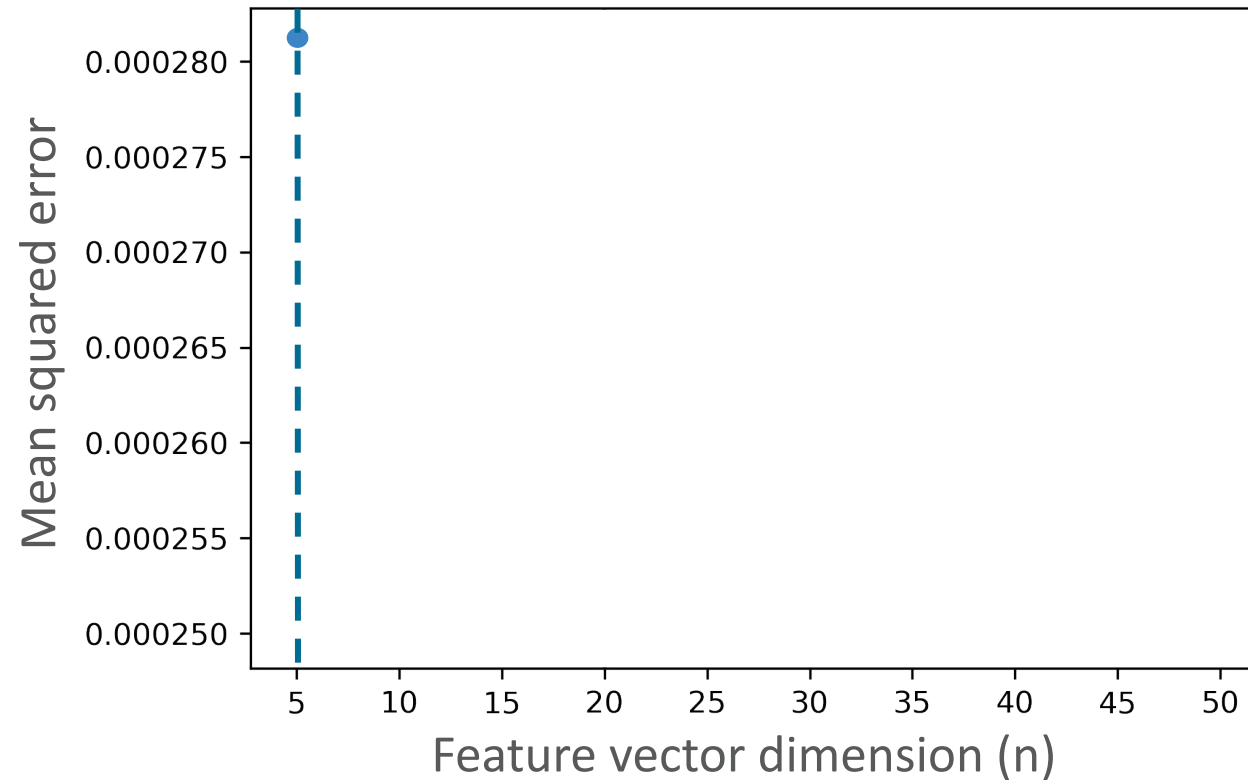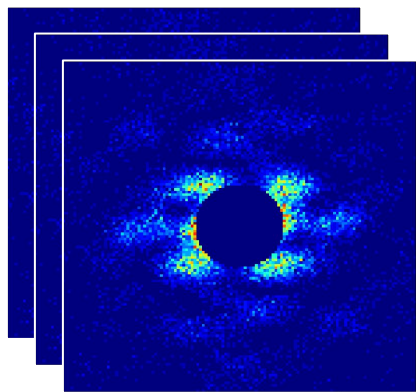


**Feature vectors**

**Reconstructed diffraction patterns**

# Feature vectors dimension (n=5)

The dimension of the feature vector has to be sufficient to faithfully reconstruct the original diffraction patterns
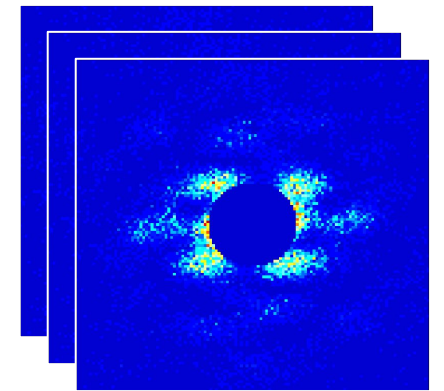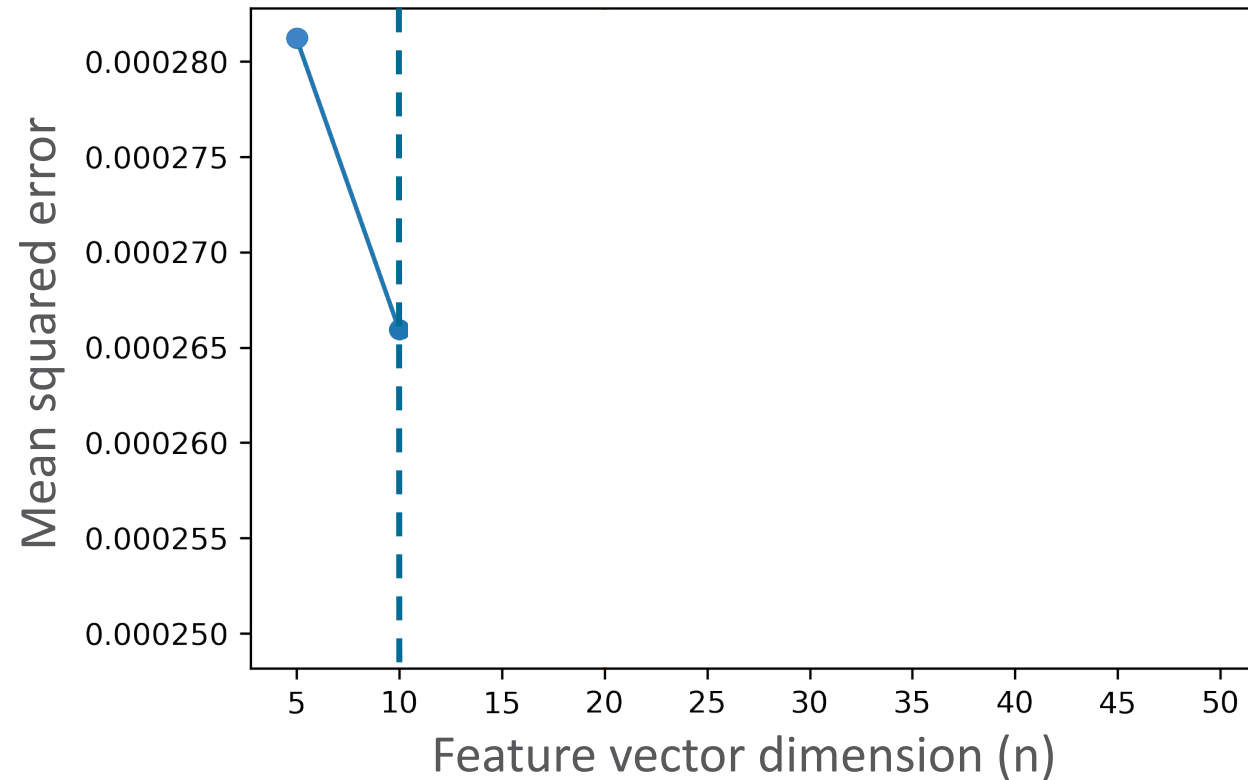


**Diffraction patterns**

**Reconstructed diffraction patterns**

# Feature vectors dimension (n=10)

The dimension of the feature vector has to be sufficient to faithfully reconstruct the original diffraction patterns
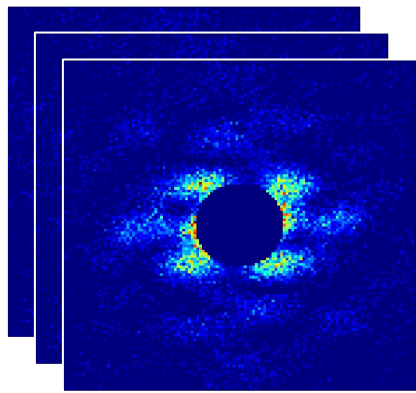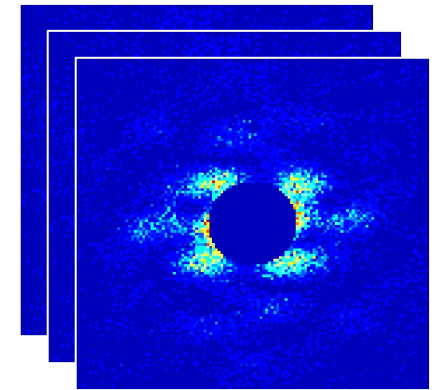


**Diffraction patterns**

**Reconstructed diffraction patterns**

# Feature vectors dimension (n=15)

The dimension of the feature vector has to be sufficient to faithfully reconstruct the original diffraction patterns
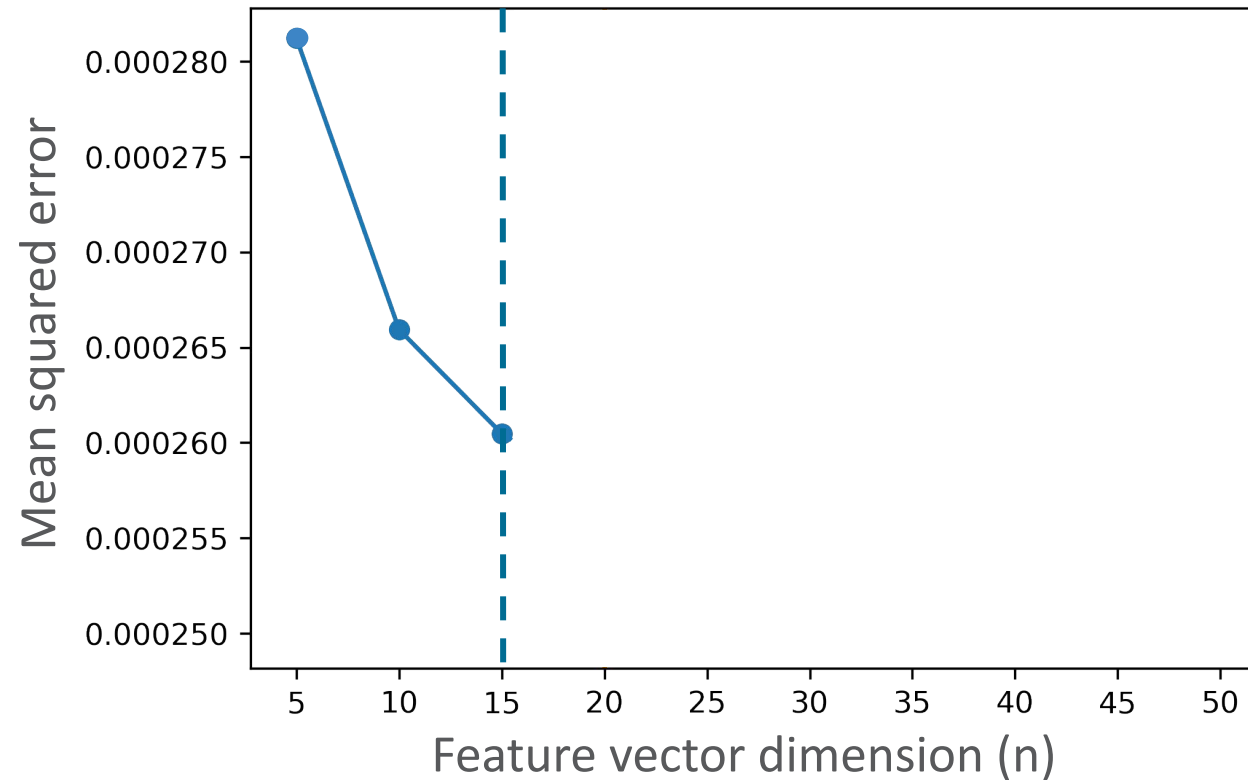


**Diffraction patterns**

**Reconstructed diffraction patterns**

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# Feature vectors dimension (n=20)

The dimension of the feature vector has to be sufficient to faithfully reconstruct the original diffraction patterns
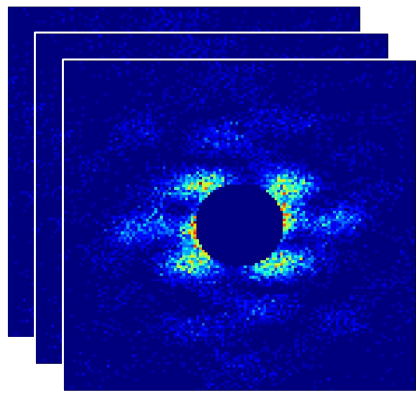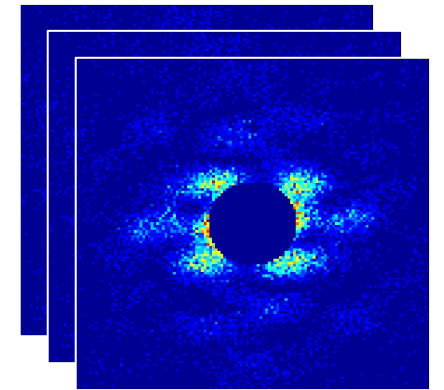


**Diffraction patterns**

**Reconstructed diffraction patterns**

# Feature vectors dimension (n=50)

The dimension of the feature vector has to be sufficient to faithfully reconstruct the original diffraction patterns
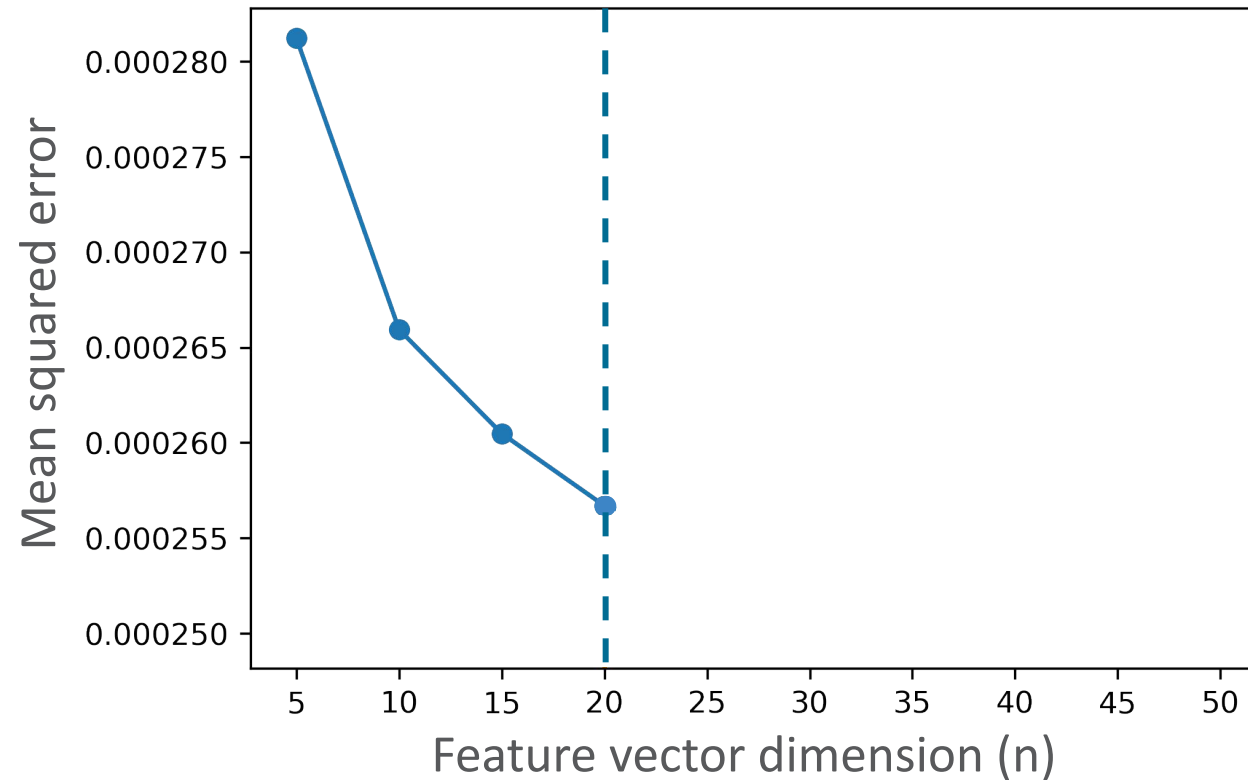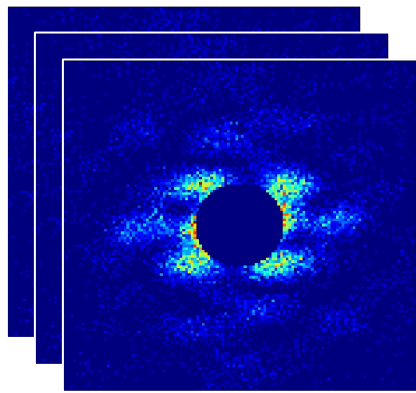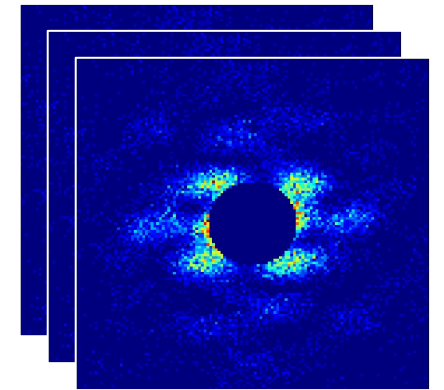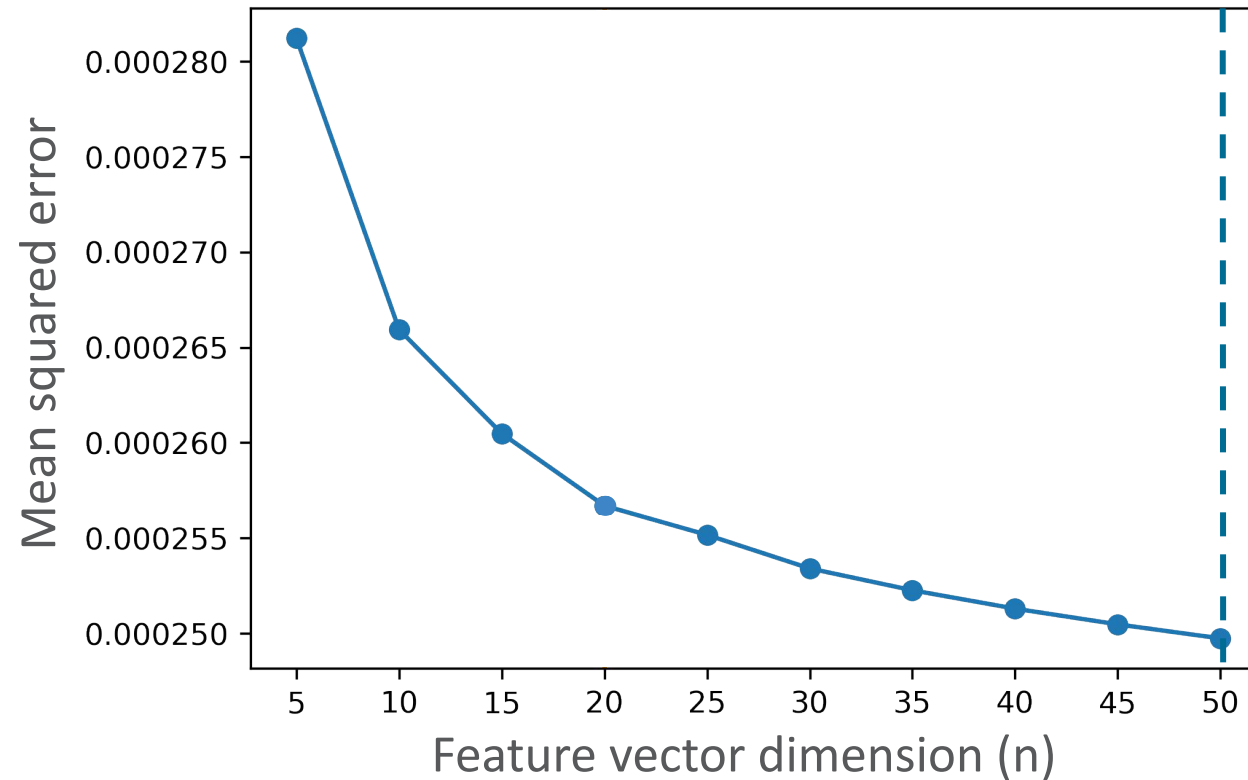


**Diffraction patterns**

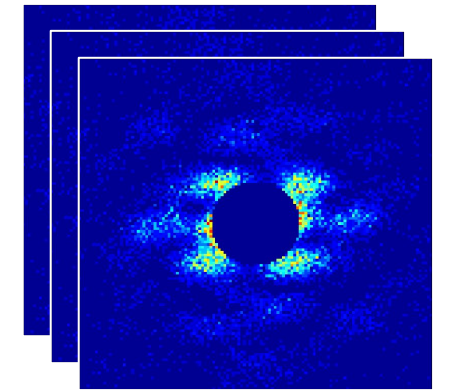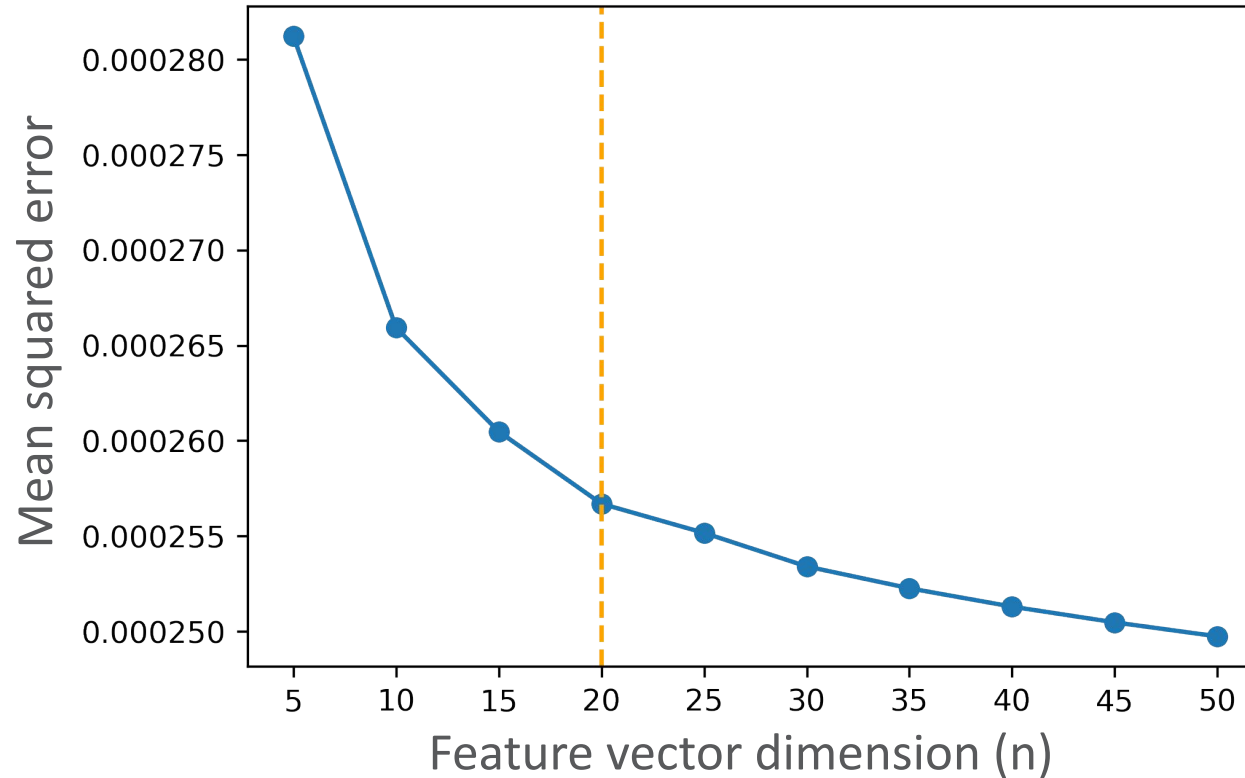**Reconstructed diffraction patterns**

# Identify the suitable feature vector dimension

Using the elbow method, our framework identifies when variance of the error and the associated gain in accuracy are not significant
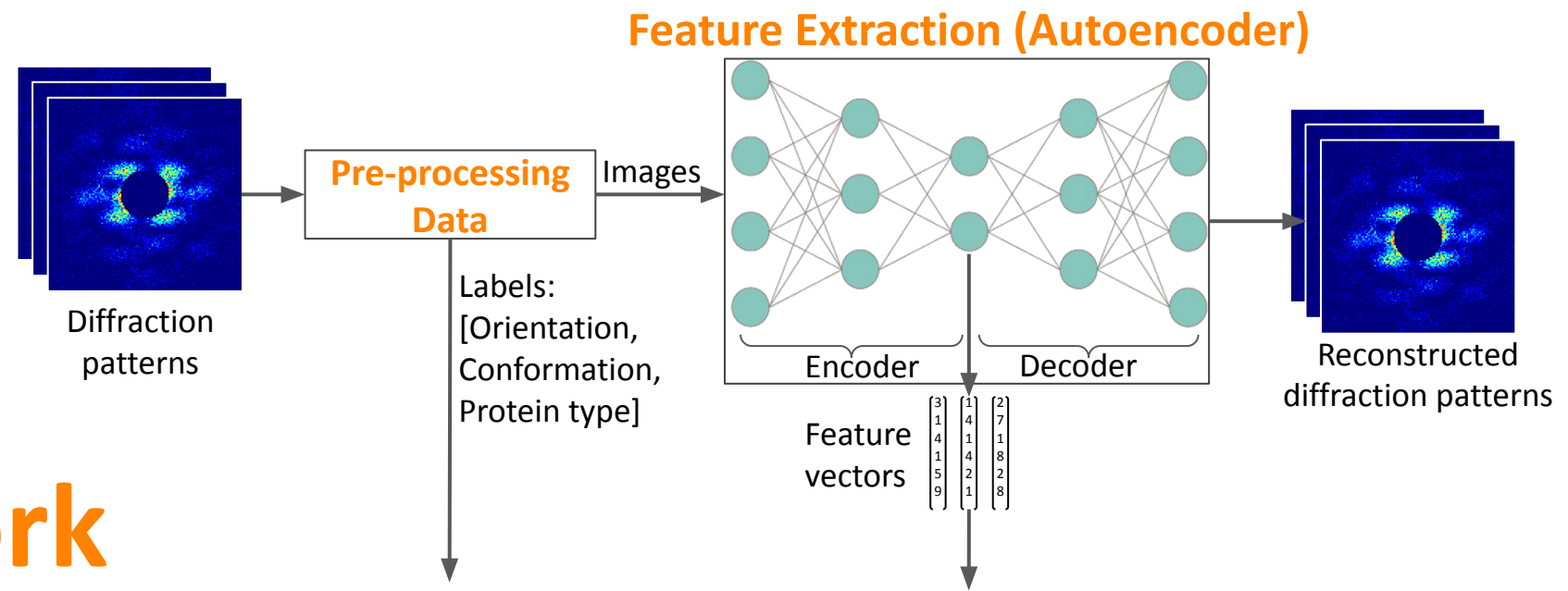
**Feature vectors**

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_{20} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_{20} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_{20} \end{bmatrix}$$
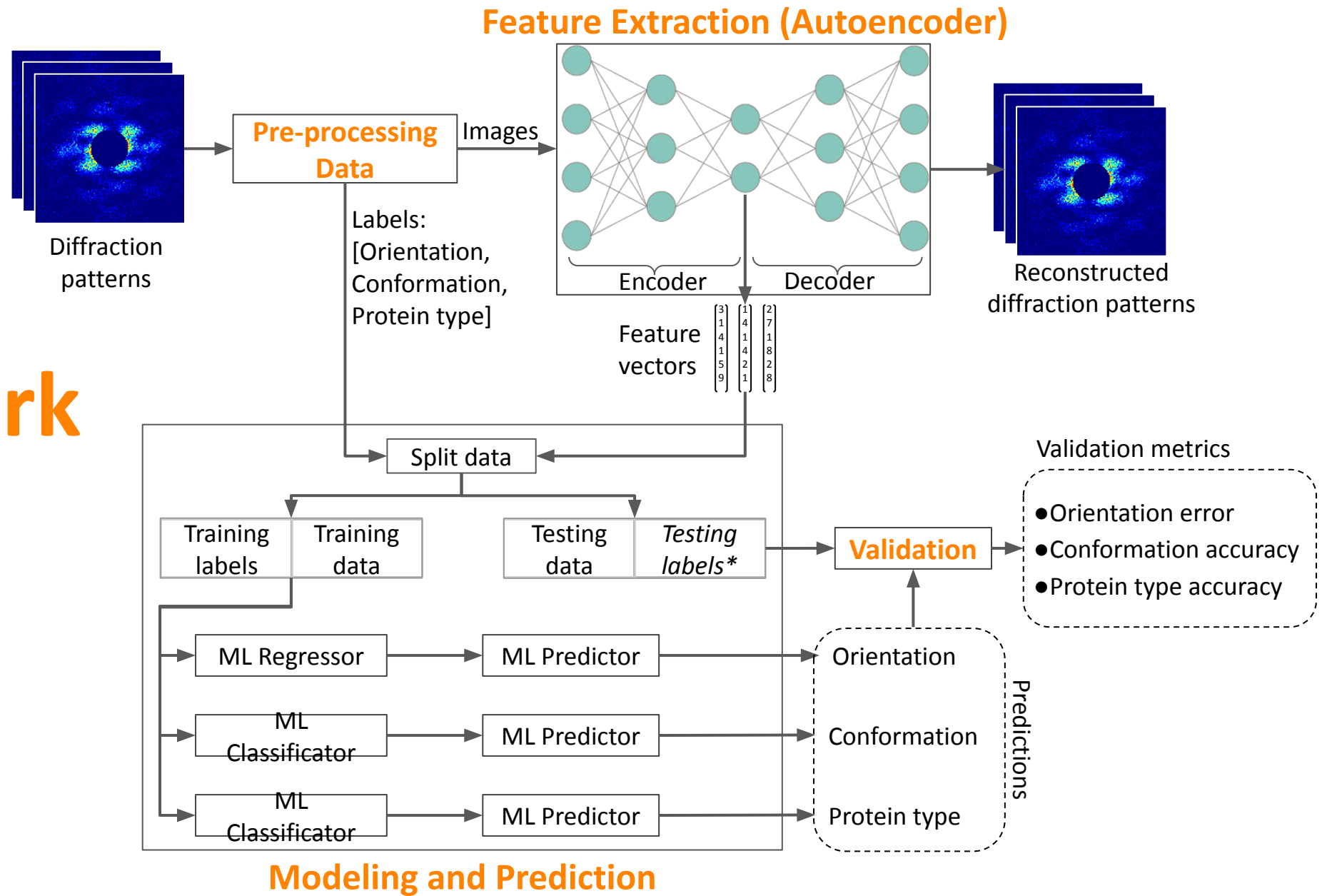
**n=20**





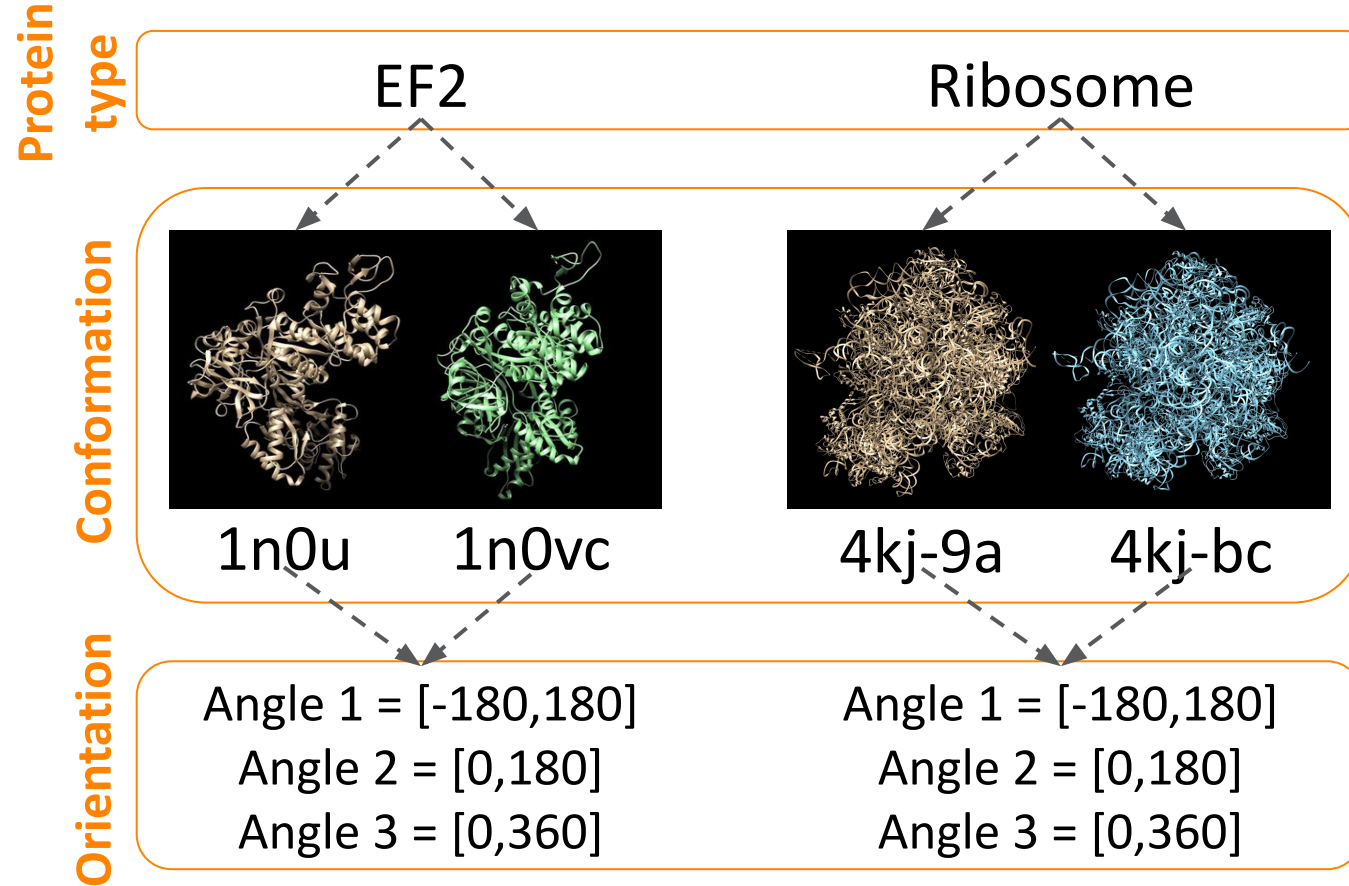**Reconstructed diffraction patterns**

XPSI Framework

# XPSI Framework

XPSI identifies structural properties (i.e., orientation, conformation, protein type)



**Feature Extraction (Autoencoder)**

Diffraction patterns

**Pre-processing Data** → Images

Labels: [Orientation, Conformation, Protein type]

Encoder | Decoder

Reconstructed diffraction patterns

Feature vectors

Split data

Training labels | Training data

Testing data | *Testing labels\**

**Validation**

Validation metrics
- Orientation error
- Conformation accuracy
- Protein type accuracy

ML Regressor → ML Predictor → Orientation

ML Classificator → ML Predictor → Conformation

ML Classificator → ML Predictor → Protein type

Predictions

**Modeling and Prediction**

THE UNIVERSITY OF TENNESSEE KNOXVILLE

# Identifying structural properties with XPSI

We demonstrate our framework's capability to identify structural properties by merging diverse datasets of diffraction patterns with multiple orientations, conformations, and protein types
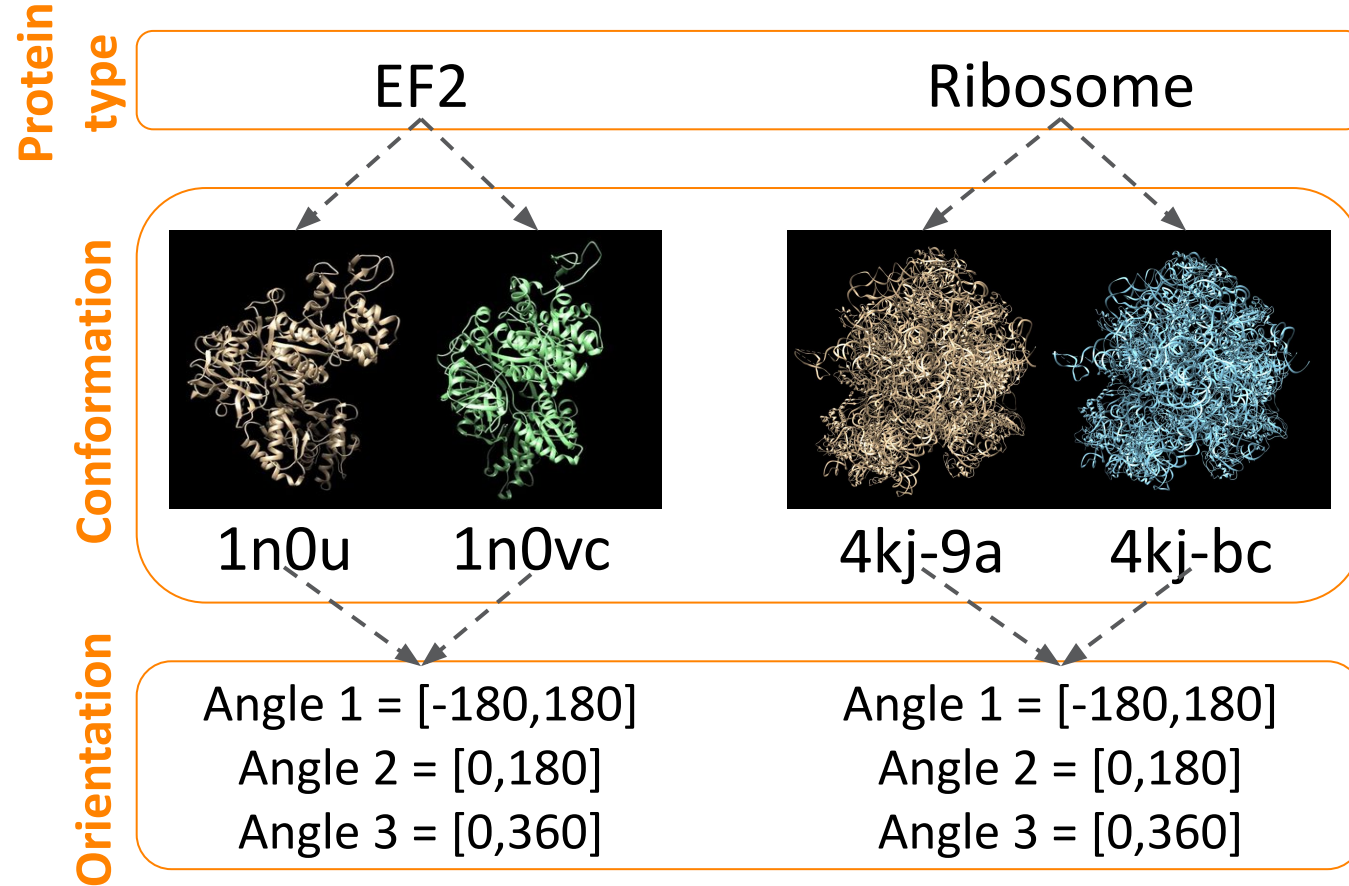
- 39,692 diffraction patterns per each conformation

THE UNIVERSITY OF TENNESSEE KNOXVILLE
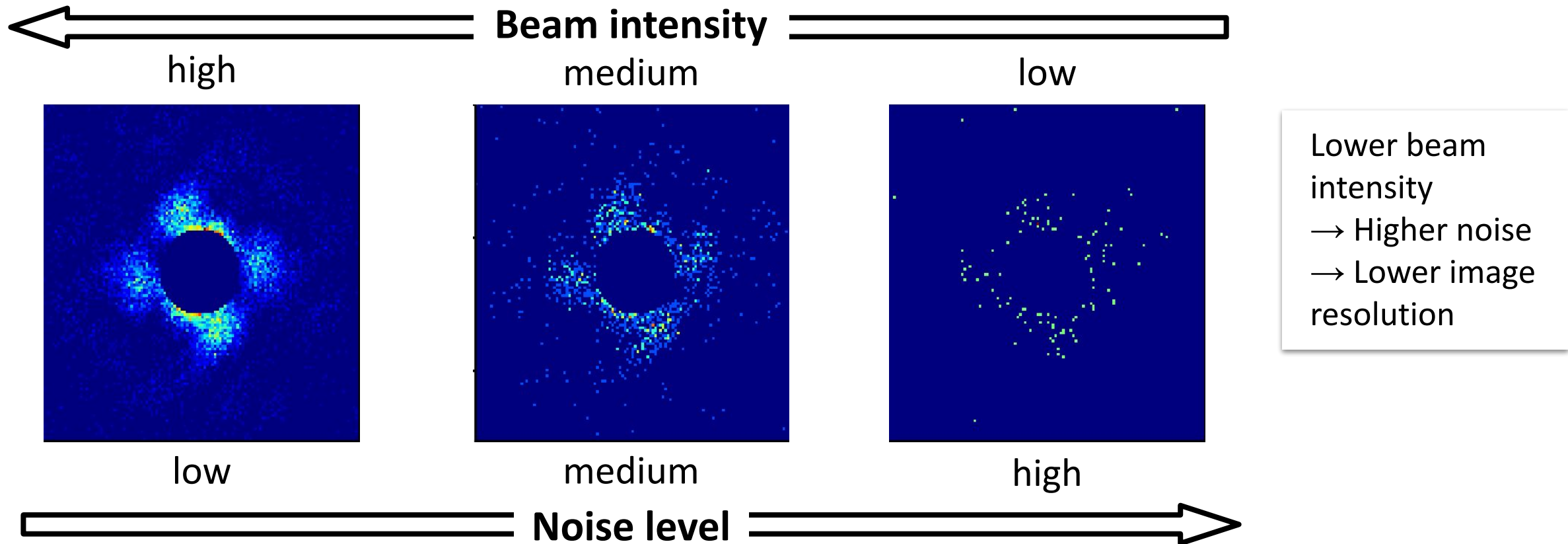
# Identifying structural properties with XPSI

We demonstrate our framework's capability to identify structural properties by merging diverse datasets of diffraction patterns with multiple orientations, conformations, and protein types
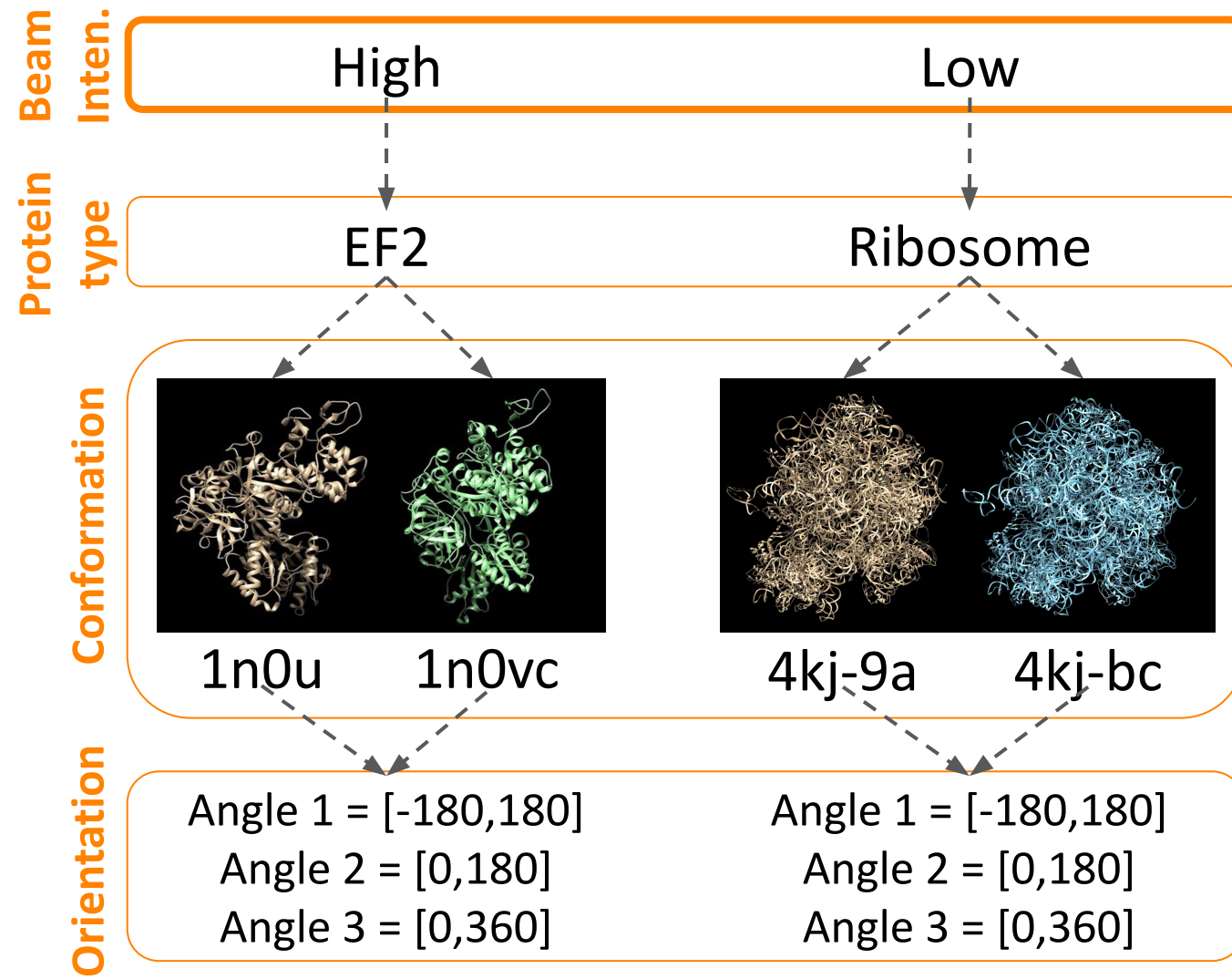
**But there is one extra challenge … Noise**



Protein type

EF2                                    Ribosome

Conformation

1n0u          1n0vc                4kj-9a          4kj-bc

Orientation

Angle 1 = [-180,180]              Angle 1 = [-180,180]
Angle 2 = [0,180]                  Angle 2 = [0,180]
Angle 3 = [0,360]                  Angle 3 = [0,360]

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# Noise in the XFEL diffraction patterns

The XFEL beam intensity is proxy for noise in the diffractions patterns (images)



**Beam intensity**
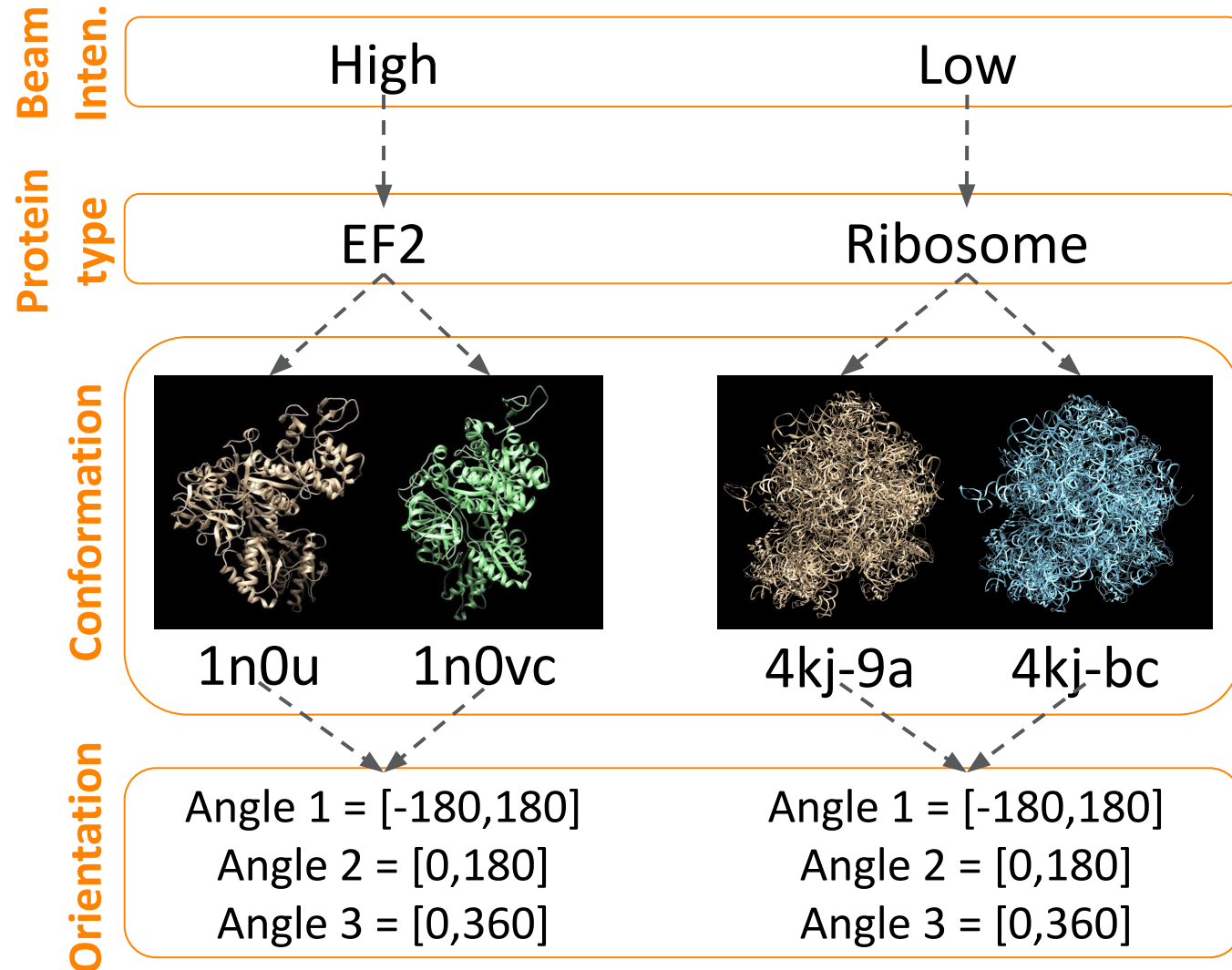
high      medium      low

low      medium      high

**Noise level**

Lower beam intensity
→ Higher noise
→ Lower image resolution

# Identifying structural properties with XPSI



**Beam Inten.**

| High | Low |

**Protein type**

| EF2 | Ribosome |

**Conformation**

1n0u    1n0vc          4kj-9a    4kj-bc

**Orientation**

Angle 1 = [-180,180]
Angle 2 = [0,180]
Angle 3 = [0,360]

Angle 1 = [-180,180]
Angle 2 = [0,180]
Angle 3 = [0,360]

# Identifying structural properties with XPSI

We quantify and validate XPSI's ability to provide accurate structural properties predictions for **diverse datasets of diffraction patterns** (multiple orientations, conformations and protein types) with **different beam intensities**

→ 10% testing data (~4000 diffraction patterns from each conformation)

**Beam Inten.**

| High | Low |

**Protein type**

| EF2 | Ribosome |

**Conformation**



| 1n0u | 1n0vc | 4kj-9a | 4kj-bc |

**Orientation**

Angle 1 = [-180,180]
Angle 2 = [0,180]
Angle 3 = [0,360]

Angle 1 = [-180,180]
Angle 2 = [0,180]
Angle 3 = [0,360]

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# Computer infrastructure

1 x 32-core Power9 node (128 GB RAM) with 1 x GPU Nvidia V100s



**Beam Inten.**
| High | Low |

**Protein type**
| EF2 | Ribosome |

**Conformation**



| 1n0u | 1n0vc | 4kj-9a | 4kj-bc |

**Orientation**

Angle 1 = [-180,180]
Angle 2 = [0,180]
Angle 3 = [0,360]

Angle 1 = [-180,180]
Angle 2 = [0,180]
Angle 3 = [0,360]

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# Orientation error

**We measure the error to predict the three angles using two metrics**

1. **Error degree**

The distance between two points on a sphere given Φ (Azimuth) and Θ (Altitude)

$$2\sqrt{\sin^2\left(\frac{\theta_2 - \theta_1}{2}\right) + \cos(\theta_1)\cos(\theta_2)\sin^2\left(\frac{\phi_2 - \phi_1}{2}\right)}$$

2. **Psi difference**

The difference between real and predicted Psi (Ψ) angle

$$\psi_{real} - \psi_{predicted}$$

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# Orientation error



The **error degree** for 90% of the testing data is below 6°

Legend:
- 50th percentile (blue dashed)
- 75th percentile (orange dashed)
- 90th percentile (green dashed)

Y-axis: Count of diffraction patterns
X-axis: Error degree °

# Orientation error



The **error degree** for 90% of the testing data is below 6º

The **psi difference** for 90% of the testing data is below 7º

THE UNIVERSITY OF TENNESSEE KNOXVILLE

# Conformation accuracy

**Accuracy:** Represents the ratio of correct predictions over the total number of cases examined

$$\frac{TP + TN}{TP + TN + FP + FN}$$

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# Conformation accuracy

**Accuracy:** Represents the ratio of correct predictions over the total number of cases examined

$$\frac{TP + TN}{TP + TN + FP + FN}$$

# Conformation accuracy

- **XPSI predicts** between 4 different **conformations with an accuracy of 90%** on average
- XPSI always predicts the conformations within the proteins (no inter-error class)

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# Protein type accuracy

- **XPSI predicts** between 2 different **protein types with an accuracy of 100%**

# XPSI remarks

We demonstrate the scientific robustness of XPSI in different challenges:

- **Identifying multiple proteins (100% of accuracy), conformations (90% of accuracy), and orientations (error degree < 6° and psi difference < 7°)**
- Differentiating between conformations (97% of accuracy) with similar, but not identical, structures of the same protein
- Identifying rotation in the diffraction patterns, even in the presence of symmetry (error degree < 10° and psi difference < 10°)

All of these capabilities are proven with different beam intensities. The lower the beam intensity the noisier the diffraction patterns, which affects the accuracy of the predictions

THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

# XPSI Jupyter notebook

We provide a Jupyter Notebook for shareability and portability of our framework

https://github.com/TauferLab/XPSI

# XFEL slice matching for 3D reconstruction

Apply our XPSI framework to XFEL slice matching for 3D reconstruction



experimental images (angles unknown)

reconstruct 3D volume with best matched angles

phase retrieval

search the best matched image to assign angle

iteration

create library by slicing

reference volume updated by iteration

reference library (angles known)

initial reference volume

Fourier space

3D structure in real space

Nakano et al., J. Synchrotron Rad. (2017). 24, 727–737

THE UNIVERSITY OF TENNESSEE KNOXVILLE

# XPSI Jupyter notebook

We provide a Jupyter Notebook for shareability and portability of our framework

https://github.com/TauferLab/XPSI

THE UNIVERSITY OF TENNESSEE KNOXVILLE