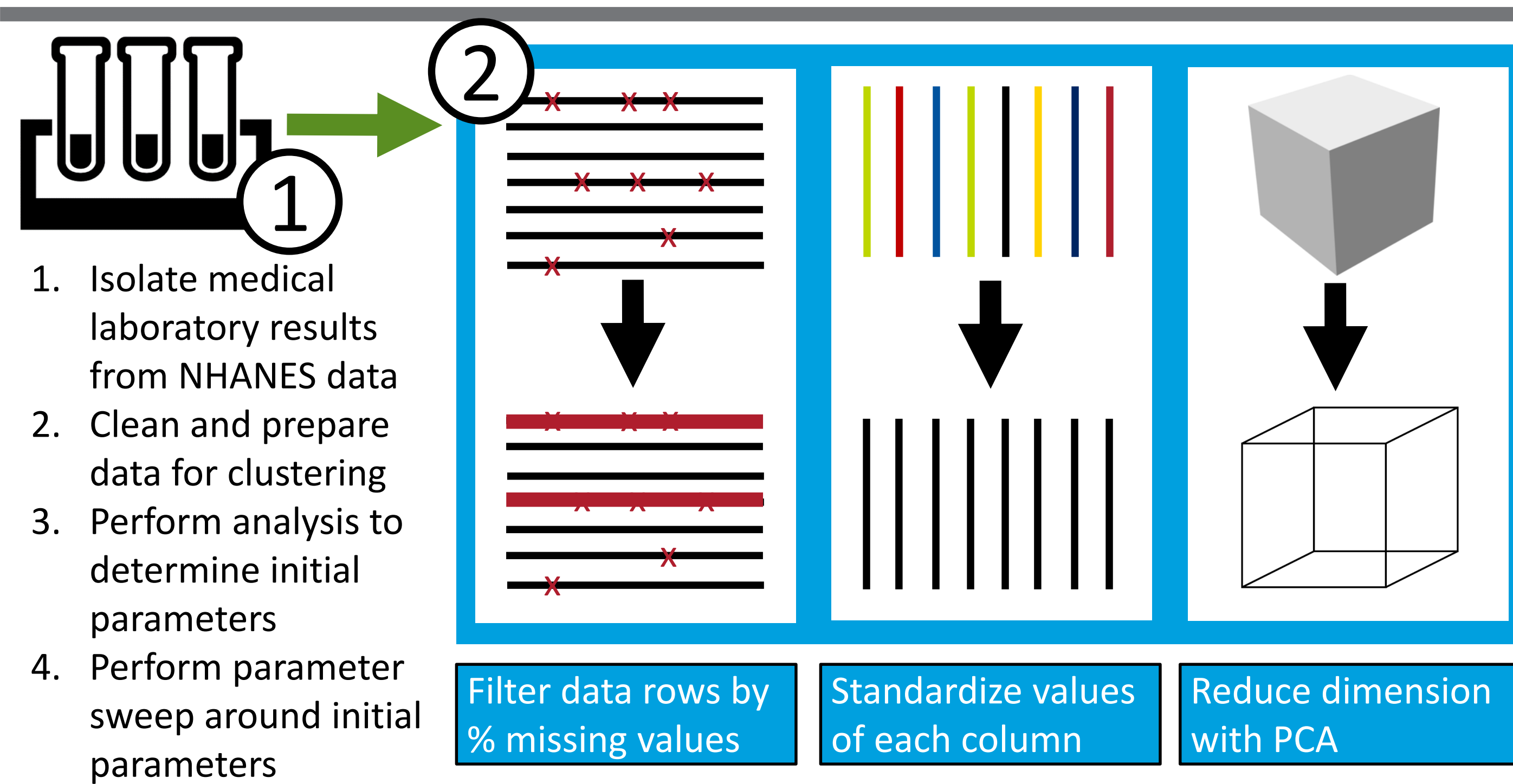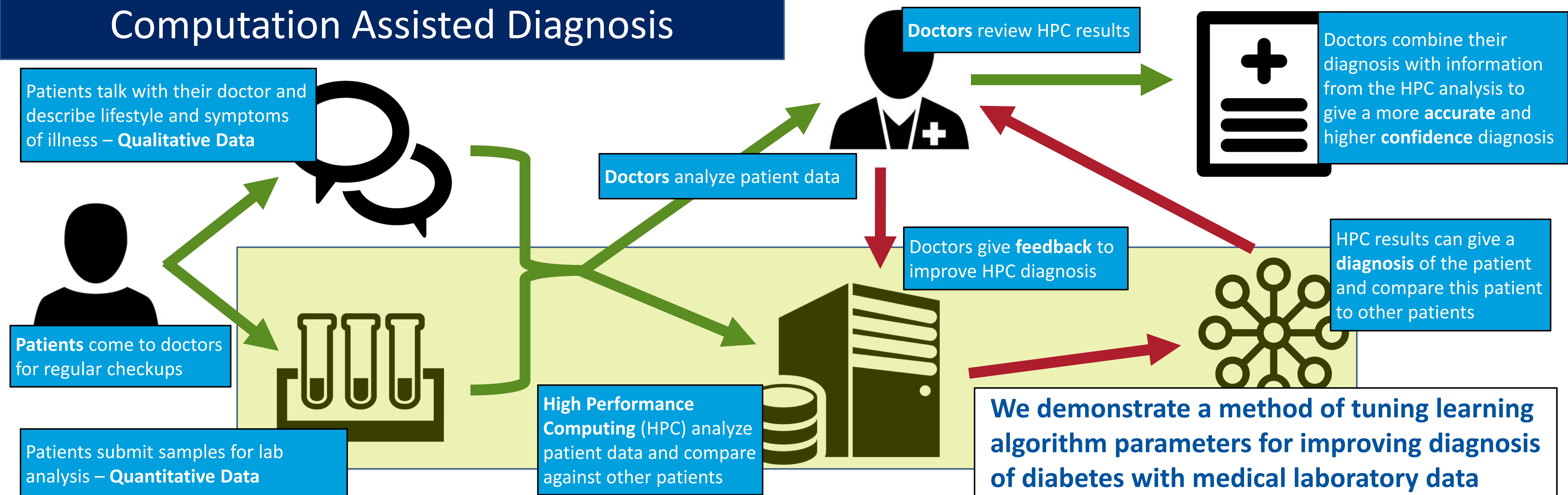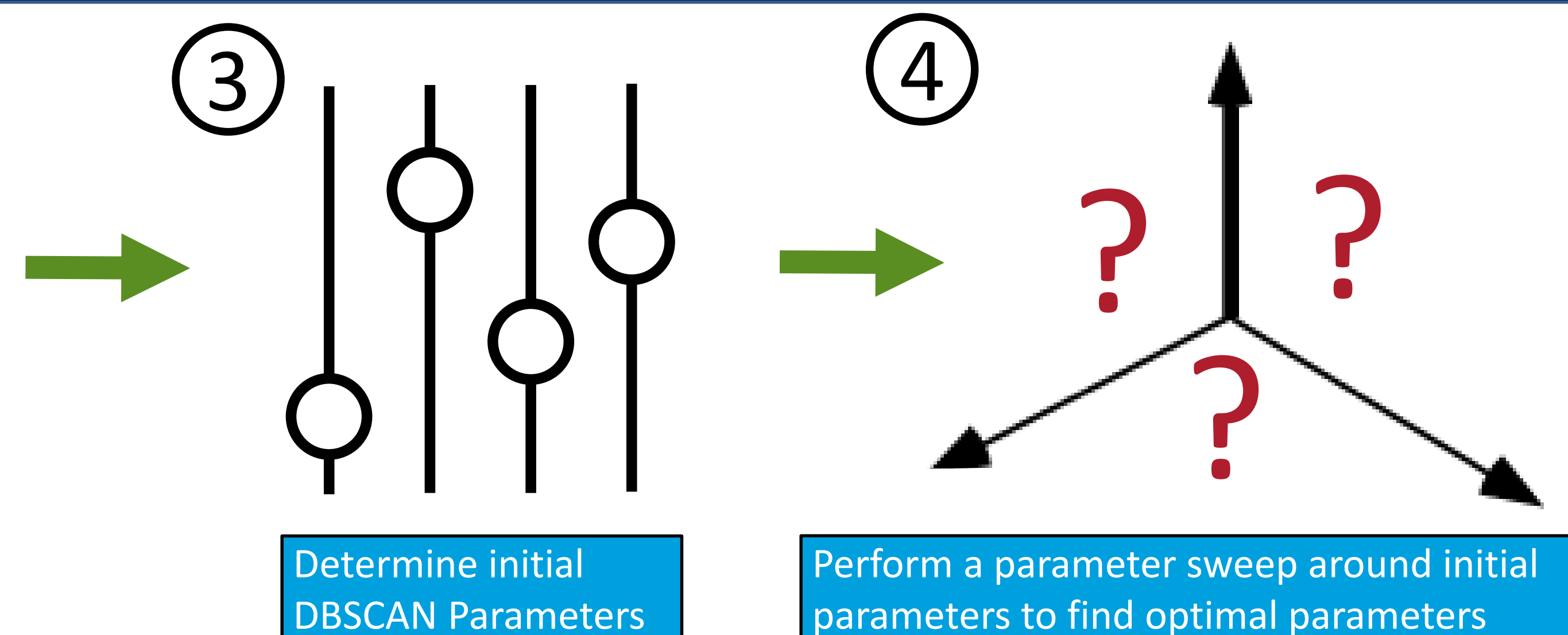# Parameter Tuning of DBSCAN for Medical Data and Diabetes Diagnosis

**Michael Wyatt, Michela Taufer**
University of Delaware: Global Computing Lab

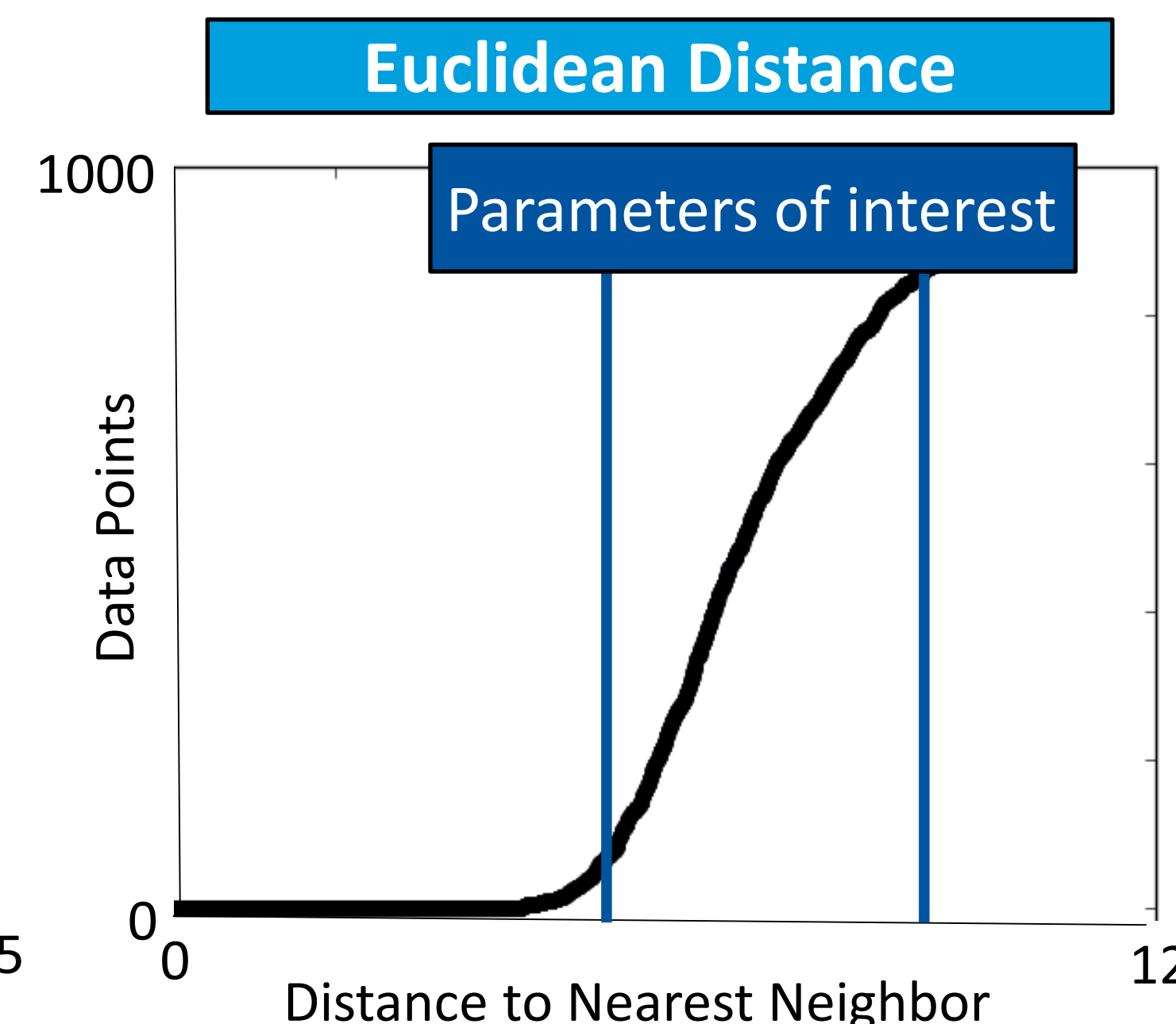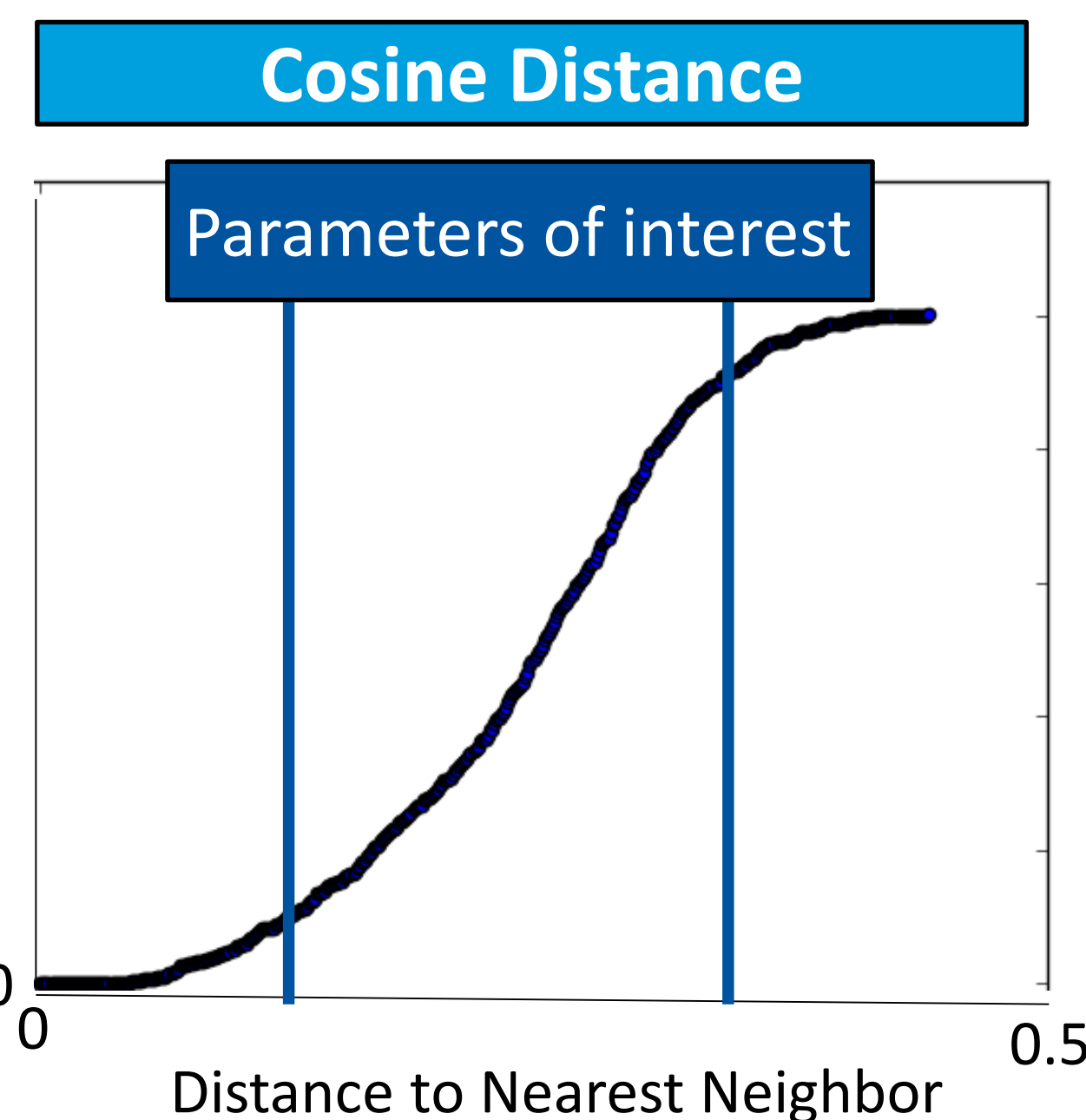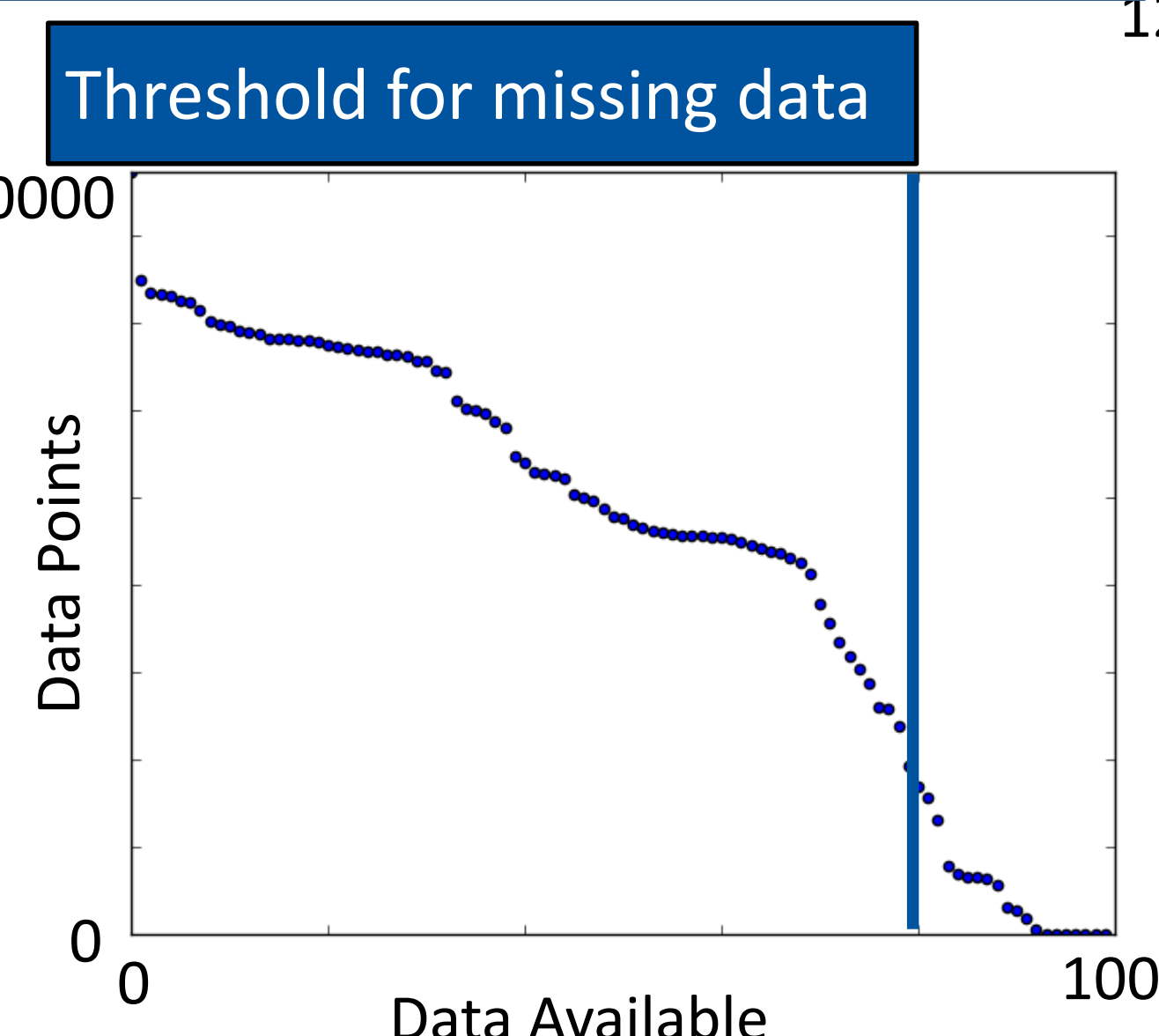## Computation Assisted Diagnosis

Patients talk with their doctor and describe lifestyle and symptoms of illness – **Qualitative Data**

**Doctors** review HPC results

Doctors combine their diagnosis with information from the HPC analysis to give a more **accurate** and higher **confidence** diagnosis

**Doctors** analyze patient data

**Patients** come to doctors for regular checkups

Doctors give **feedback** to improve HPC diagnosis

HPC results can give a **diagnosis** of the patient and compare this patient to other patients

Patients submit samples for lab analysis – **Quantitative Data**

**High Performance Computing** (HPC) analyze patient data and compare against other patients

**We demonstrate a method of tuning learning algorithm parameters for improving diagnosis of diabetes with medical laboratory data**

1. Isolate medical laboratory results from NHANES data
2. Clean and prepare data for clustering
3. Perform analysis to determine initial parameters
4. Perform parameter sweep around initial parameters

Filter data rows by % missing values

Standardize values of each column

Reduce dimension with PCA

## Parameter Tuning Workflow

③ Determine initial DBSCAN Parameters

④ Perform a parameter sweep around initial parameters to find optimal parameters

## Parameter Tuning Results

- Missing data **trade-off**
  - No Missing Values → **Small Dataset**
  - Many Missing Values → **Bad Clusters**

  **Must find a balance between missing values and dataset size**
- Patients with > 80% lab data available
  - 16,627 patients (over 1/5 original data)
  - Produces higher quality clusters

- DBSCAN parameters affect cluster quality
  - **Epsilon**: Neighborhood to search for neighbors
  - **Min_pts**: minimum neighbors to be in a cluster

  Parameters define cluster **density**
- Distance metrics also affect cluster quality: **Euclidean** vs. **Cosine**

- Utilizing nearest neighbor analysis, we can determine the range of epsilon values which should be tested
- We cluster data with several **epsilon** and **min_pts** values around the identified optimal values
- We measure the quality of each clustering by **percentage of points clustered** and **information gained** by each clustering:

$$score = \frac{Diebetes\ Patients\ Clustered}{Total\ Diabetes\ Patients} * \frac{Informaion\ Gain}{Max\ Information\ gain}$$

- We identify an optimal parameter setting:
  - **Cosine distance, Epsilon: 0.15, Min_pts: 3**

**Threshold for missing data**

Data Points — Data Available (0 to 100), 0 to 90000

**Cosine Distance** — Parameters of interest — Data Points (0 to 1200), Distance to Nearest Neighbor (0 to 0.5)

**Euclidean Distance** — Parameters of interest — Data Points (0 to 1000), Distance to Nearest Neighbor (0 to 12)

**Cosine Distance** — Min points (2 to 6), Epsilon (0.15 to 0.35), 14%

**Euclidean Distance** — Min points (2 to 6), Epsilon (5 to 9), 2.5%