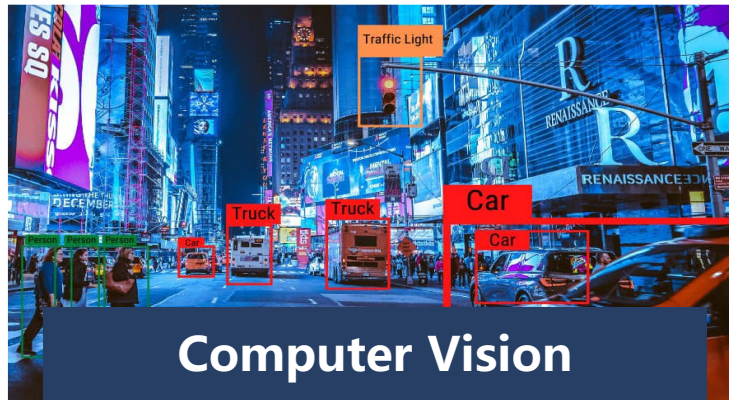


## **Introduction to Data Science (Part 1)**

Stan Z. Li

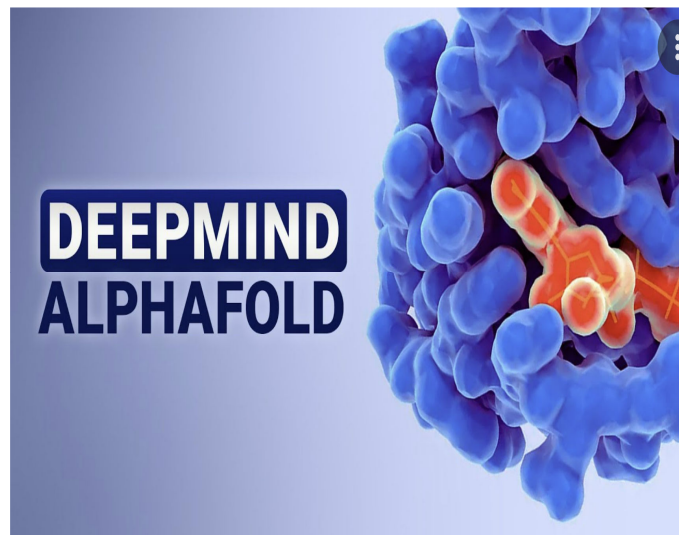
# AI Research (Main-Stream Tasks)



# AI for Sciences



**Optimization/  
Decision Making**



**Protein Folding  
and Design**



**Drug Design  
and Synthesis**

# Examples of AI Modeling Problems

- Face Recognition
- Speech
- Natural Language Processing
- Robot Sensing and Control
- Whether forecasting
- Protein Computing (Structure Prediction and Sequence Design)
- Mass-spectrometric Data Analysis
- Single-Cell Clustering and Hierarchy/Lineage
- Drug Design (Large and Small Molecules)
- .....

# 西湖大学AI 三大研究版块

## AI 基础研究

- 数据科学基础
- 深度学习方法
- 序列结构建模

## AI 核心应用

- 计算机视觉
- 语音语言处理
- 机器人学

## AI 学科交叉

- 生命科学
- 生物医药
- 其他 学科

# AI and Machine Learning

## ARTIFICIAL INTELLIGENCE

Any technique that enables computers to mimic human behavior



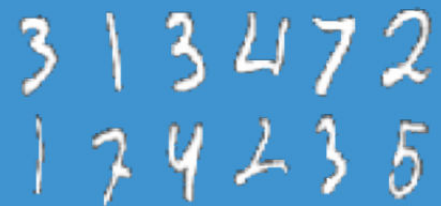
## MACHINE LEARNING

Ability to learn without explicitly being programmed



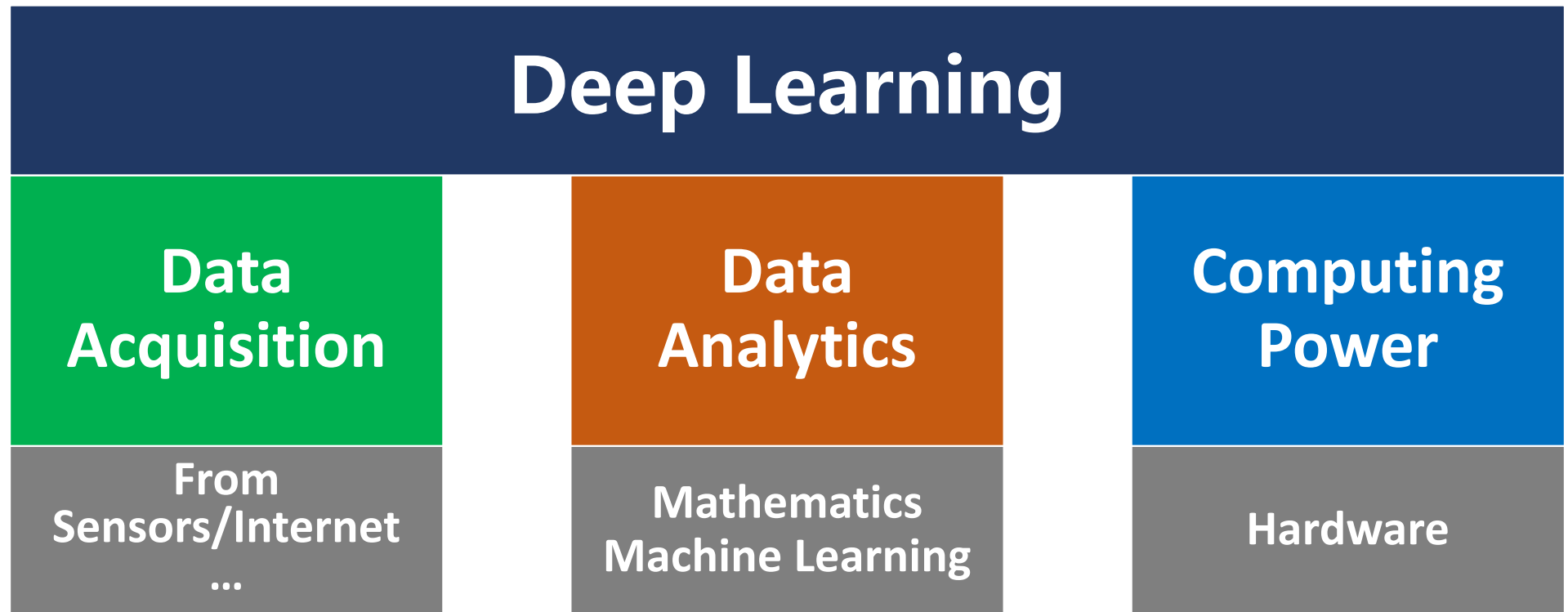
## DEEP LEARNING

Extract patterns from data using neural networks



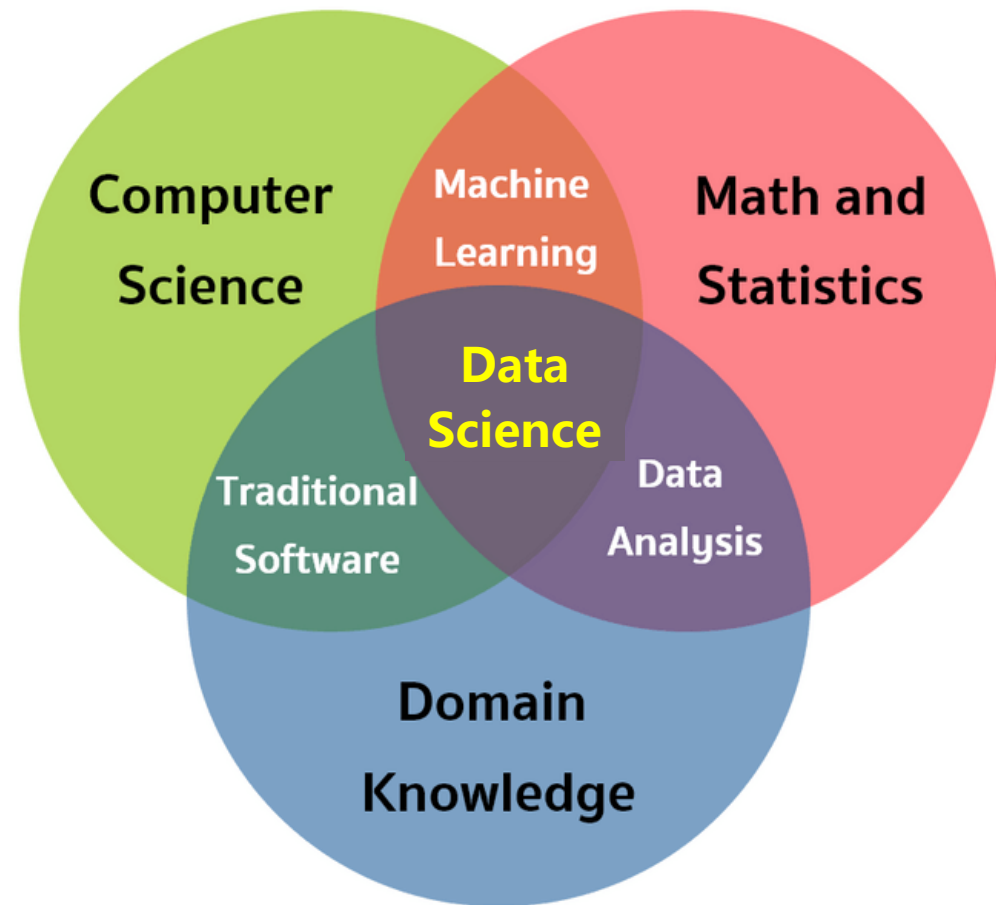
Picture Courtesy of MIT 6.S191

# Contemporary AI and Deep Learning



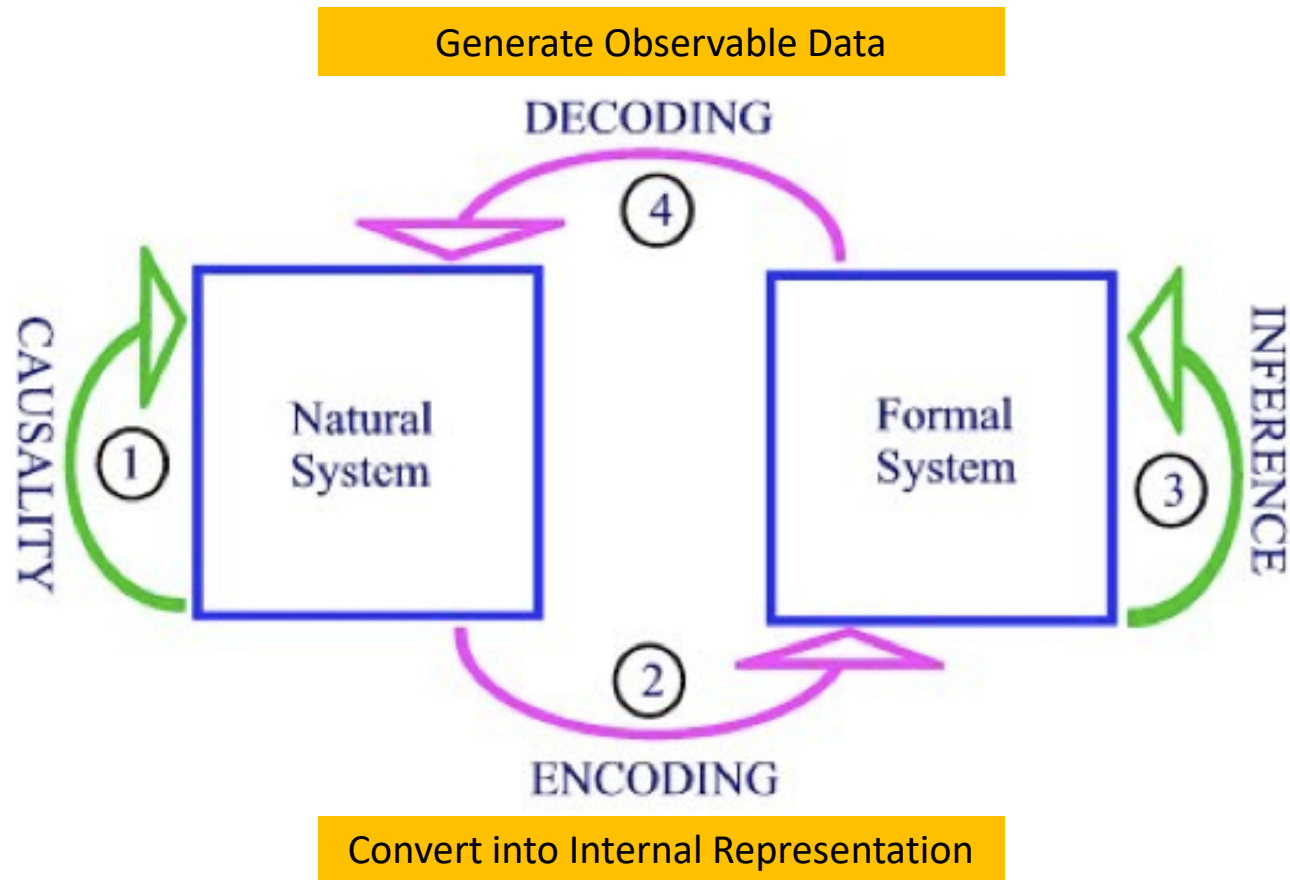


# Data Science





# Scientific Modeling



# Different Types of Models

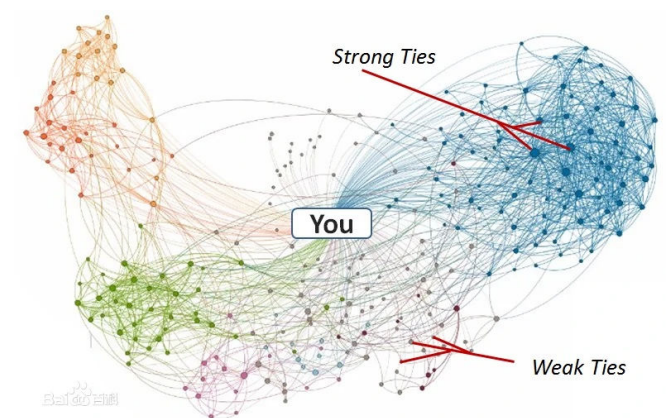
- **Conceptual models:** to better understand
- **Operational models:** to operationalize
- **Mathematical models:** to quantify
- **Computational models:** to simulate  
(and to encode / decode)
- **Graphical models:** to visualize the subject

# Outline

1. AI – How it has been evolving
2. **High-Dimensional Data Analysis**
3. Modeling by Deep Learning and Neural Nets

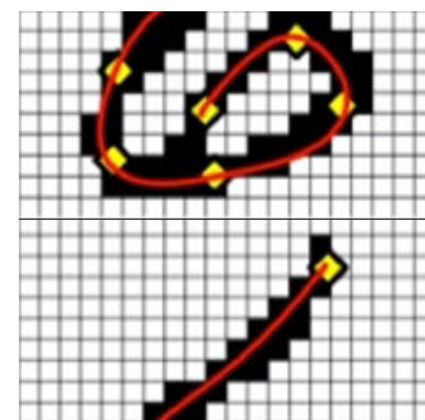
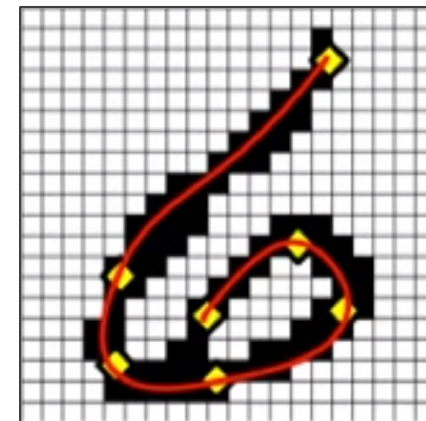
# High-Dimensional Data

- Images, Videos, Text, Audio,
- Web pages, Social Networks
- Molecular Structures
- DNA Sequences
- Protein Sequence-Structures



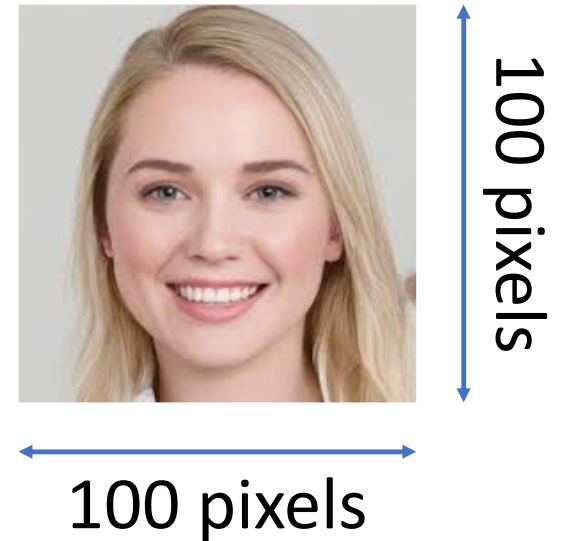
# Handwritten Digit images

- Image size  $20 \times 20 = 400$
- Pixel values in  $\{0,1\}$
- Image Space  $\mathcal{S} = \{0,1\}^{400}$
- $\#\mathcal{S} = 2.58 \times 10^{120}$
- Only a tiny portion of  $\mathcal{S}$  is of digits
- The digit pattern lives in a low dim subspace (manifold)



# Face Image Data

- Image size  $100 \times 100 = 10^4$  pixels
- RGB image size  $3 \times 10^4$  pixels
- **Dimensionality =  $3 \times 10^4$**
- Pixel values in  $\{0, \dots, 255\}$
- #Possibility =  $256^{30,000} \cong \text{infinity}$
- Only a tiny portion is of faces
- **Face pattern lives in low dim subspace**



# Manifold Assumption

**High-Dimensional Data:** Images, Web pages, Gene sequences, ....

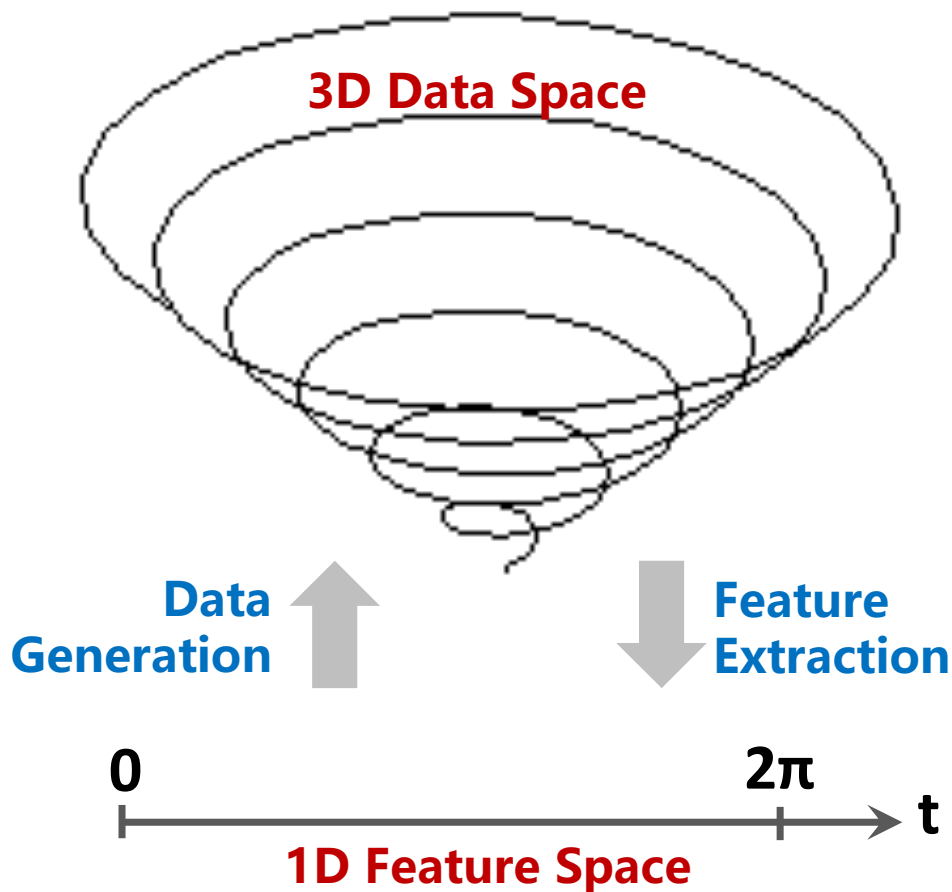
**Dimension Reduction into Coordinate System of a Lower Dim**

- For representation learning (feature extraction)
- For data visualization – in 2D or 3D

**Manifold Assumption: an interesting pattern in high dimensional data resides on a low dimensional manifold**



# Manifold in Hi-D Data Space: 1D Curve in 3D Space



**Conical Helix:**

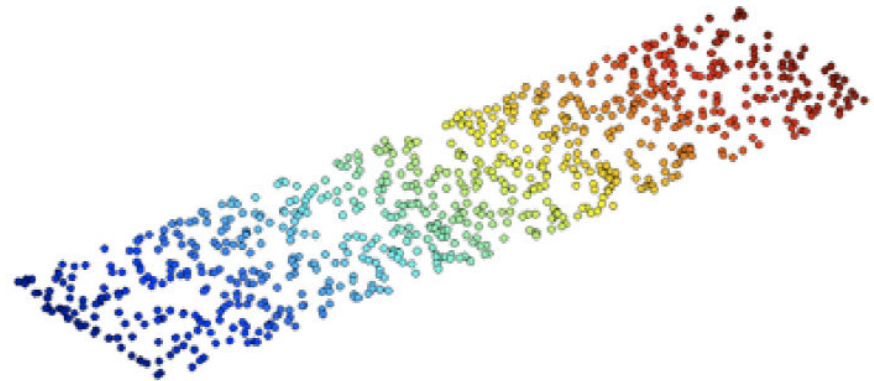
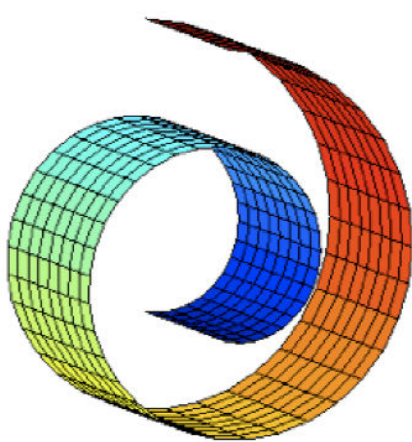
$$x=t*\cos(6t), y=t*\sin(6t), z=t$$

$$0 \leq t \leq 2\pi$$

**1D line segment**

**Latent variable  $t$**

# 2D Manifold in 3D Space

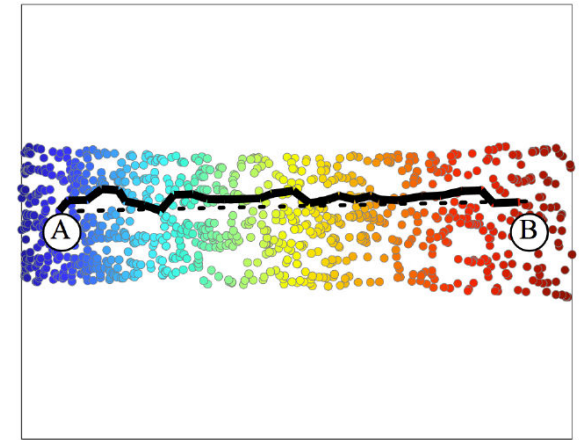
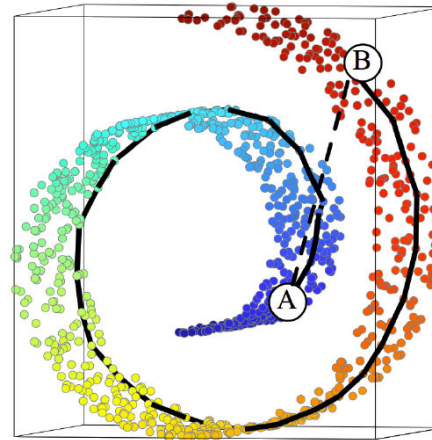
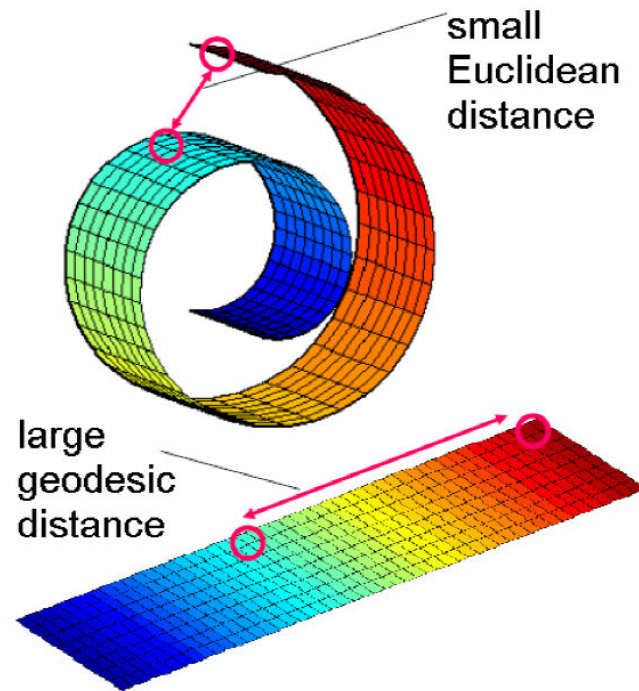


## Swiss Roll:

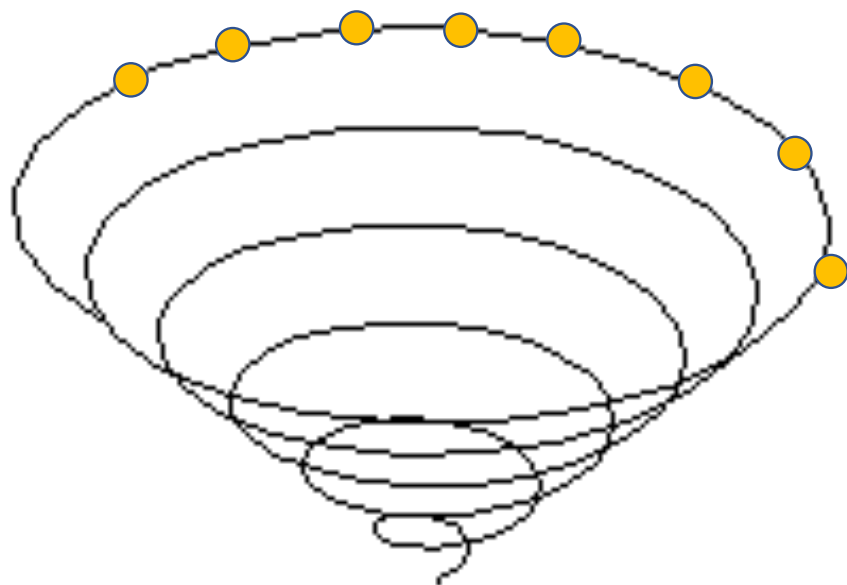
$$x = \varphi \cos(\varphi), y = \varphi \sin(\varphi), z = \psi$$
$$1.5\pi \leq \varphi \leq 4.5\pi, 0 \leq \psi \leq 10$$

**Manifold: 2D rectangle**  
generated by two latent  
variables  $\varphi, \psi$

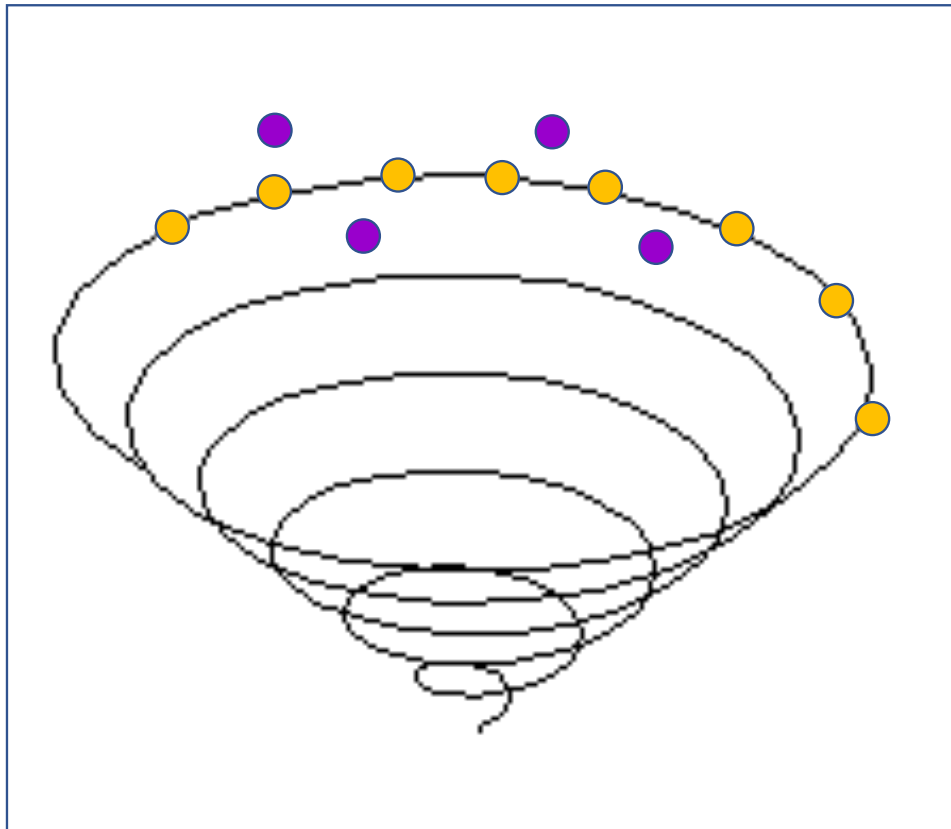
# Geodesic Distance on Manifolds



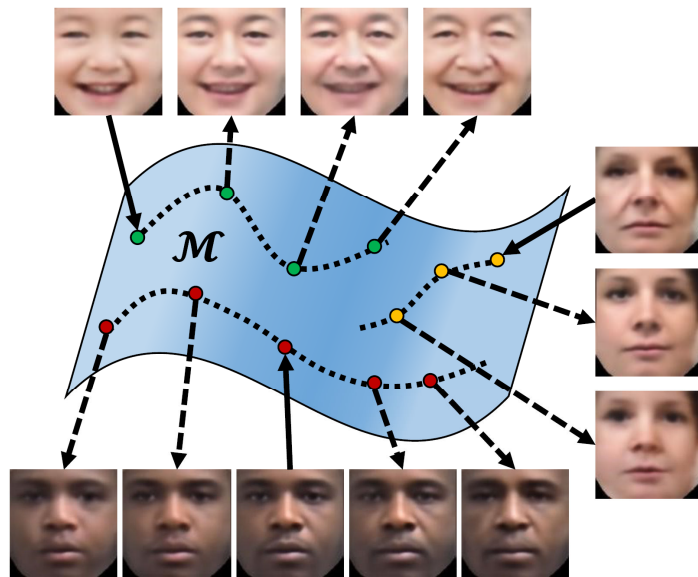
# Samples in Data Space



# Samples Close to the Face Manifold



# 2D Surface in 3D Space



特征提取变换  
Feature Extraction

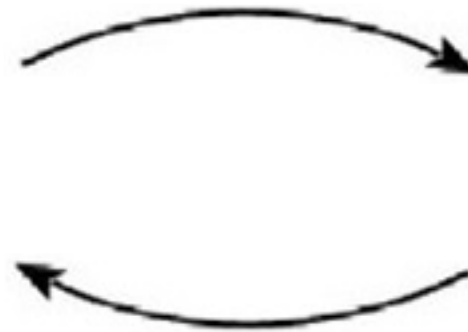


Image Generation  
图像生成变换

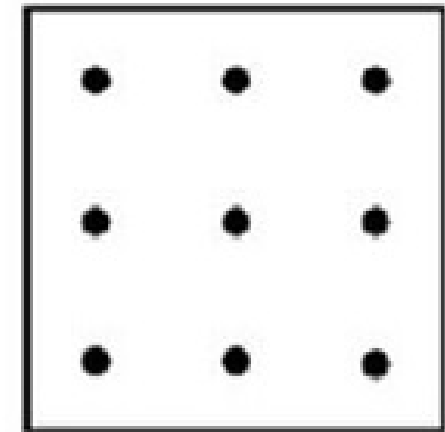


Image Space:  $10^4$ -D

Feature Space:  $10^2$ -D

**Q & A**



# Thanks

