## **Pursuing the Nature of Intelligence**

## Professor Yi Ma School of Computing and Data Science The University of Hong Kong





#### Seek a scientific and theoretical foundation for Intelligence:

- what to learn?
- how to learn?
- why correct?

"What I cannot create, I do not understand."

-- Richard Feynman





#### **Evolution of Life and Intelligence**

Evolution of the Universe is **Physics** at work



#### Evolution of Life is Intelligence at work



"Just as the constant increase of entropy is the basic law of the universe, so it is the basic law of life to be ever more highly structured and to struggle against entropy." -- Vaclay Havel

#### **Evolution of Life and Intelligence: From DNA to Brain**

From the first DNA to the emergence of life with Brain: **3.6 Billion Years** From the first Brain to the explosion of lives in the Cambrian period: **50 Million Years** 



A Brief History of Intelligence, Max Bennett, 2023

#### **Evolution of Life and Intelligence: From Species to Individuals**

#### Emergence and evolution of life are mechanisms of intelligence at work!

Life depends on intelligence to continuously acquire more knowledge to better predict the world.

**Phylogenetic Intelligence:** DNA inheritance, random mutation, and natural selection

**Ontogenetic Intelligence:** Individual memory, perception & feedback, and error correction.



3.7B years ago Life begins 500M years ago Cambrian period

400M years ago fish 360M years ago amphibian 250M years ago reptile

200M years ago bird and mammal

310T years ago neanderthal

#### **Evolution of Intelligence: From Societal to Artificial Intelligence**

#### Emergence and evolution of life are mechanisms of intelligence at work!

Life depends on intelligence to continuously acquire more knowledge to better predict the world.

**Societal Intelligence:** Languages and texts, empirical knowledge, trial and error

#### Artificial Intelligence: Scientific facts, theorize, hypothesis testing & falsification.



310T years ago neanderthal

Tools and group hunting

70T years ago Societal intelligence

Languages Information sharing

Knowledge accumulation

3500 BC

written language

600-300 BC **Artificial Intelligence** 

14-18<sup>th</sup> Century Renaissance

Abstraction, formal logic, and mathematics

Science

Hypothesis Testing

The 1940s Machine Intelligence

Computing machines

### The Origin of Machine Intelligence (the magic era!)

1940s, people started to make machines imitate intelligence (of animals).

- 1948, Cybernetics & System Theory, Nobert Wiener
- 1943, Artificial Neural Networks, Warren McCulloch and Walter Pitts
- 1948, Information Theory, Claude Shannon
- 1944, Game Theory, John von Neumann
- 1940's, Turing Machine and Turing Test, Alan Turing





#### **Artificial Neurons and Neural Networks: Learn from Nature**

#### Dendrite Axon terminal Soma (cell body) x Outputs Mvelin sheat Output points = synapses Myelinated axon trunk Inputs Input points = synapses $x_0 = +1$ $x_1$ $w_{k0} = b_k$ $x_2$ $v_k$ $\varphi(\cdot)$ $x_3$ $y_k$ $w_{L2}$ $\tilde{u}_{km}$ $x_m$

Golgi and Cajal 1888 (1901 Nobel Prize)

#### Hubel and Wiesel 1959 (1981 Nobel Prize)



Warren McCulloch & Walter Pitts 1948

#### Fukushima 1980 & LeCun 1989 (Turing Award)

### History of Machine Intelligence (Artificial Neural Networks)



#### Figure courtesy of Professor Rene Vidal

### **Modern Evolution of Deep Neural Networks**



### Why Must Turn Blackbox to Whitebox?

٠

•

٠

- Modern AI systems all based on empirically designed deep networks (alchemy?)
- Blackbox is difficult to explain, impossible to guarantee, costly to improve, ...

#### It is high time to develop a principled approach!



### What to Learn?

#### The fundamental reason why intelligence exists and evolves: The world is not entirely random yet, and it is still largely predictable.

**Intelligence and Science** learn what is predictable from sensed data of external world (so every animal is Newton and has learned an accurate "world model.")







### What to Learn?

#### The fundamental reason why intelligence exists and evolves: The world is not entirely random yet, and it is still largely predictable.

Mathematically, all predictable information is encoded as a distribution p(x) of **low-dimensional** supports in observed **high-dimensional** data space.



#### This is the only "inductive bias" necessary!

High-dimensional Data Analysis with Low-Dimensional Models, Wright and Ma, 2022

### What to Learn: Low-dimensionality

Important properties of low-dimensional structures: Completion

Bayes inference: 
$$y = \mathcal{P}_{\Omega}(x) \rightarrow \widehat{x} \sim p(x \mid y)$$

**Conditional sampling** 





Image completion







Corruptions

Corruptions

### What to Learn: Low-dimensionality

Important properties of low-dimensional structures: Denoising Empirical Bayes:  $y = x + \sigma n \implies \hat{x} \sim \mathbb{E}(x \mid y) = y + \sigma^2 \nabla \log p(x)$ 

Tweedie's formula

Noisy observations



Natural image denoising (or generation)





### What to Learn?

Important properties of low-dimensional structures: Error Correction

**Robust Bayes inference**:  $y = x + e \rightarrow \hat{x} \sim \operatorname{argmin} ||x||_1 + ||e||_1$ 

Exploiting sparsity prior

Corrupted observations





## What to Learn: Seeking Parsimony

#### The main objective of learning:

Identify a distribution with low-dimensional structures from sensed data x and transform them to a compact and structured representation z.



## How to Learn: Measure of Parsimony

For distributions with **low-dim** structures, **entropy** (or "volume") of the underlying data distribution should be very small (or "zero"):

(Discrete) Entropy: 
$$H(X) = \sum_{x \in X} -p(x) \log p(x)$$

Differential Entropy: 
$$h(\mathbf{x}) = \int -p(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}$$





### **How to Learn: Compress**

A fundamental and unifying mechanism to learn low-dim structures: compress to reduce **entropy** of the observed (noisy) data distribution.

minimize 
$$h(\mathbf{x}) = \int -p(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}$$



### How to Learn: Compress via Denoising

Theorem [Diffusion] Consider the diffusion process:

$$x_t = x_o + tn, \quad n \sim \mathcal{N}(0, I).$$

Under natural technical assumptions, the **entropy of the process increases**:

 $\frac{d}{dt}h(\boldsymbol{x}_t) > 0, \qquad \forall t > 0.$ 



Learning Deep Representations of Data Distributions, Ma+PBWY, 2025

### How to Learn: Compress via Denoising

Theorem [Denoising] Consider the inverse denoising process:

$$\widehat{\boldsymbol{x}}_{t-s} = \mathbb{E}[\boldsymbol{x}_{t-s} \mid \boldsymbol{x}_t] = \boldsymbol{x}_t + st\nabla \log p(\boldsymbol{x}_t)$$

Under natural technical assumptions, the entropy of the process decreases:



Learning Deep Representations of Data Distributions, Ma+PBWY, 2025

### How to Learn: Empirical Approaches



## **How to Learn: Empirical Approaches**

Empirically designed networks to realize the denoising operator:

 $\tilde{z}^{l+1} = \tilde{z}^l + \beta \nabla \log p(\tilde{z}^l)$  how to realize?



Diffusion Transformer (DiT)

Image or video generation



Latent Diffusion Transformer

DiT Block with adaLN-Zero

### How to Learn: An Analytical Case

Analytically derived operation to realize the denoising operator:

 $\tilde{z}^{l+1} = \tilde{z}^l + \beta \nabla \log p(\tilde{z}^l)$  how to realize?

If we approximate a general distribution  $p(\tilde{z})$  with a **mixture of subspaces or low-dim Gaussians,** e.g. PCA, ICA, GPCA, Sparse Coding [W+Ma, 2022],

$$\tilde{\boldsymbol{z}}^l \sim \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\boldsymbol{0}, \boldsymbol{U}_k \boldsymbol{U}_k^{\mathsf{T}}), \qquad \boldsymbol{U}_k \in O(D, d),$$

then

$$\tilde{\boldsymbol{z}}^{l+1} \propto \frac{1}{K} \sum_{k=1}^{K} \operatorname{softmax} \{ \alpha || \boldsymbol{U}_{i}^{\mathsf{T}} \tilde{\boldsymbol{z}}^{l} ||^{2} \}_{k} \boldsymbol{U}_{k} \boldsymbol{U}_{k}^{\mathsf{T}} \tilde{\boldsymbol{z}}^{l} .$$



## How to Learn: Measure of Information Gain

differential entropy

How to explicitly represent a distribution with a low-dimensional support? A key idea: **lossy encoding and decoding** 

construct a finite codebook by packing the support of the distribution with  $\epsilon$ -balls.

 $h(\mathbf{x}) = -\infty$   $R(\mathbf{x}, \epsilon) = \min_{\mathbb{E}||\hat{\mathbf{x}} - \mathbf{x}|| \le \epsilon} h(\hat{\mathbf{x}}) - h(\hat{\mathbf{x}} | \mathbf{x})$   $\overset{2\epsilon}{\longleftrightarrow}$   $\overset{2\epsilon}{\longleftrightarrow}$   $S_{1}'$ 

Segmentation of Mixed Data via Lossy Coding and Compression [Ma+DHW, TPAMI2007].

rate distortion

## How to Learn: Organize Information

How to make the resulting representation the most informative?

compress what is similar; contrast what is dissimilar.



Mathematically: maximize information gain or reduce coding rate:  $\max \Delta R(z, \epsilon) = R(z, \epsilon) - R^{c}(z, \epsilon)$ 

ReduNet: A White-box Deep Network from the Principle of MCR<sup>2</sup> [YCY+Ma, JMLR2022].

#### How to Learn: Maximal Coding Rate Reduction (ReduNet)

When  $Z = f(X; \theta)$  is a mixture of K Gaussians with  $\Pi_k$  encodes membership of samples in the kth class.

$$\begin{array}{l} \text{Maximal Coding Rate Reduction (MCR2)} \\ \max_{f} \Delta R(Z \mid \Pi) = \frac{1}{2} \log \det \left( I + \frac{d}{\varepsilon^2} \cdot \frac{ZZ^{\top}}{m} \right) - \underbrace{\sum_{k=1}^{K} \frac{\text{tr}(\Pi_k)}{2m} \log \det \left( I + \frac{d}{\varepsilon^2} \cdot \frac{Z\Pi_k Z^{\top}}{\text{tr}(\Pi_k)} \right)}_{R(Z)} \\ \end{array}$$

Unrolled optimization: gradient ascent on the  $MCR^2$  objective,



ReduNet: A White-box Deep Network from the Principle of MCR<sup>2</sup> [YCY+Ma, JMLR2022].

### How to Learn: Maximal Coding Rate Reduction (ReduNet)

#### Benign local and global optimization landscape of MCR<sup>2</sup>

Maximal Coding Rate Reduction (MCR <sup>2</sup> )	Regularized Maximal Coding Rate Reduction (MCR <sup>2</sup> )
$\max_{f} \Delta R(\mathbf{Z} \mid \mathbf{\Pi}) = R(\mathbf{Z}) - R^{c}(\mathbf{Z} \mid \mathbf{\Pi})$	$\max_{f} F(\mathbf{Z}) := \Delta R_{\lambda}(\mathbf{Z} \mid \mathbf{\Pi}) = R(\mathbf{Z}) - R^{c}(\mathbf{Z} \mid \mathbf{\Pi}) - \lambda \cdot \ \mathbf{Z}\ _{F}^{2}$
Theorem [YCY+M NeurIPS2020].	Theorem [WLY+M 2024].
The global optimal solution $Z^* = [Z_1^*, Z_2^*,, Z_K^*]$ of MCR <sup>2</sup> satisfies: • Subspaces of different classes are orthogonal to each other, $(Z_i^*)^T Z_j^* = 0$ for $i \neq j$ ;	Every critical point $\{Z : \nabla F(Z) = 0\}$ is either a local maximizer or a strict saddle point: each local maximizer corresponds to a feature representation that consists of a family of orthogonal subspaces.
• Each subspace achieves its maximal dimension, i.e., rank( $\mathbb{Z}_k^{-}$ ) = $a_k$ , and and $a = \sum_{k=1}^{k} a_k$ .	global maximizer





### How to Learn: Sparse Rate Reduction (SRR)

Maximize the difference between the coding rate of all features and the coding rate of features within each subspace, and promote sparsity:

$$\max_{\mathbf{Z}} \operatorname{SRR}(\mathbf{Z} | \mathbf{U}_{[K]}) = R(\mathbf{Z}) - R^{c}(\mathbf{Z} | \mathbf{U}_{[K]}) - \lambda \|\mathbf{Z}\|_{1}$$

 $-\lambda ||\mathbf{Z}||_1$ : measure how sparse all features are.



## How to Learn: DNN to Realize Iterative Optimization

Optimize Sparse Rate Reduction via gradient descent (GD)

 $\min L_{SRR}(\boldsymbol{Z}) = -SRR(\boldsymbol{Z} \mid \boldsymbol{U}_{[K]})$ 





gradient descent  $\mathbf{Z}^{\ell+1} = f^{\ell}(\mathbf{Z}^{\ell}) \approx \Pr[\mathbf{Z}^{\ell} - \eta \cdot \nabla L_{srr}(\mathbf{Z}^{\ell})]$ 

Sparse Rate Reduction (SRR):  $\min_{f} L_{srr}(Z \mid U_{[K]}) = \underbrace{R^{c}(Z \mid U_{[K]})}_{compression} + \frac{\lambda ||Z||_{1} - R(Z)}{sparsification}$ 

Design the  $\ell$ -th layer  $f^{\ell} = f_2^{\ell} \circ f_1^{\ell}$  via an Alternating Minimization Scheme:

- Compression Step:  $\mathbf{Z}^{\ell+1/2} = f_1^{\ell}(\mathbf{Z}^{\ell}) \approx \mathbf{Z}^{\ell} \eta \cdot \nabla R^c(\mathbf{Z}^{\ell}; U_{[K]});$
- Sparsification Step:  $\mathbf{Z}^{\ell+1} = f_2^{\ell}(\mathbf{Z}^{\ell+1/2}) \approx \operatorname{Prox}_{\lambda \|\cdot\|_1} \left[ \mathbf{Z}^{\ell+1/2} \eta \cdot \operatorname{grad}(\mathbf{Z}^{\ell+1/2}) \right].$

## How to Learn: Interpretation of Each Layer



## How to Learn: DNN to Realize Iterative Optimization

Optimize Sparse Rate Reduction via iterative gradient descent (GD)

 $f: X = Z^0 \xrightarrow{f^1} Z^1 \xrightarrow{f^2} Z^2 \longrightarrow \cdots \xrightarrow{f^L} Z^L = Z$  $\min L_{SRR}(\mathbf{Z}) = -SRR(\mathbf{Z} \mid \mathbf{U}_{[K]})$ Ζ.  $\mathbf{L}(\mathbf{Z})$ Transformer block - L. . . →Z<sup>ℓ+1</sup> →Z<sup>ℓ</sup> Transformer block -  $(\ell + 1)$  $\mathbf{Z}^{\ell} \mathbf{Z}^{\ell+1}$ Transformer block -  $\ell$ . . .  $\mathbf{Z}^{\ell+1} = \left( \mathbf{Z}^{\ell} - \eta \cdot \nabla L(\mathbf{Z}) \right)$ Transformer block - 1 fľ X

Each layer of a deep network (e.g., Transformer) realizes a GD operator.

### How to Learn: Interpretation of the Whole Network

- Forward encoding: given fixed subspaces and dictionaries  $(U_{[K]}^{\ell}, D^{\ell})_{\ell \in [L]}$ , each layer performs compression and sparsification on representations.
- Backward learning the "codebook": backpropagation to learn subspaces and dictionaries  $(U_{[K]}^{\ell}, D^{\ell})_{\ell \in [L]}$  from data.



## How to Learn: Better Semantic Interpretability

Not only mathematically fully explainable, but also semantically more interpretable!



#### How to Learn: Better Networks from First Principles

No more trial and error to design better network architectures:

- Explainable [NeurIPS 2023]
- Scalable [NeurIPS 2024] ۲

Accuracy on ImageNet-1K

Top-1

- More efficient [ICLR 2025] ۲
- More compact [CPAL 2025] ٠



#### How to Learn: Better Networks from First Principles

No more trial and error to design better network architectures: SimDINO: Simplifying DINO via Coding Rate Regularization [ICML2025]



#### Classification (ImageNet)

Method	Model	Epochs	<i>k</i> -NN	I Linear	
DINO	ViT-B	100	72.9	76.3	
SimDINO	ViT-B	100	74.9	77.3	
DINO	ViT-L	100	_	_	
SimDINO	ViT-L	100	75.6	77.4	
DINOv2	ViT-B	100	76.0	77.2	
SimDINOv2	ViT-B	100	78.1	79.7	
DINOv2	ViT-L	100	80.8	82.0	
SimDINOv2	ViT-L	100	81.1	82.4	
SwAV	ViT-S	800	66.3	73.5	
MoCov3	ViT-B	300	_	76.7	



#### Segmentation (Microsoft COCO)

		Detection ↑		Segmentation $\uparrow$			
Method	Model	$AP_{50}$	$AP_{75}$	AP	$AP_{50}$	$AP_{75}$	AP
SimDINO	ViT-L/16	5.4	1.9	2.4	4.5	1.4	1.9
SimDINO	ViT-B/16	5.2	2.0	2.5	4.7	1.5	2.0
DINO	ViT-B/16	3.9	1.5	1.8	3.1	1.0	1.4
DINO	ViT-B/8	5.1	2.3	2.5	4.1	1.3	1.8

	Hyperparameter	SimDINOv2	DINOv2	SimDINO	DINO	
	Patch size		16			
	Register tokens	4		0		
Madal	Pos-embedding anti-alias	True		False		
Model	Init layer scale	0.1	1e-5	-		
	Drop path rate	0.3		0.1		
	Weight normalize last layer	removed	True	removed	True	
	Output prototypes K	removed	65536	removed	65536	
	Init EMA momentum	0.9	0.992	0.99	5	
	Centering temperature	removed	0.07	removed	0.07	
	Warm-up temperature	removed	0.04	removed	0.04	
Dinalina	Warm-up temperature epochs	removed	30	removed	30	
Pipeine	iBOT sample prob.	0.5		-		
	iBOT mask ratio	0.1-0.	.5	-		
	iBOT head tying	False		-		
	Koleo loss weight	removed	0.1	-		
	Global crops scale	0.4 - 1				
	Local crops scale	0.05 - 0.4				
Data	Local crops number	10				
	Global crops size	224				
	Local crops size	96				
	Batch size	128x	8	64x8		
	Epochs	100				
	Warm-up epochs	10				
	Freeze last layer epochs	removed	1	removed	1	
Optim.	Learning rate	0.004		0.002		
	Layerwise lr decay	0.9		-		
	Weight decay	0.04				
	Weight decay end		0.4			
	Gradient clip	3.0		0.3		

## How to Learn: Summary with a Comparison

No more trial and error to design better network architectures:



## Why Correct? (Consistency)

**Bi-directional encoding and decoding (e.g., compression and generation)** 

a post



(2024 Nobel Prize Physics)

## Why Correct? (Consistency)

#### **Bi-directional encoding and decoding (e.g., compression and generation)**

## Masked autoencoding with a **whitebox architecture** [PBWY+**Ma**, ICLR2024]

#### $\mathcal{P}_{\Omega}(\mathbf{x}) \rightarrow \widehat{\mathbf{x}} \approx \mathbf{x}$





Image completion with 75% patches/pixels masked



Models: CRATE-MAE-Base versus ViT-MAE-Base [PBWYM, ICLR2024]

## **Towards Autonomous Intelligence (AI 2.0)**

How to self-learn a more consistent representation, continuously?



In nature, all intelligent systems learn from closed-loop feedback! (Cybernetics)



## **Towards AI 2.0: Close the Loop via Minimax Game**

#### Closed-loop systems learnt via minimax game do not forget catastrophically!

Incremental Learning via Closed-Loop [TDWLY+Ma, ICLR 2023]

Unsupervised Learning of Structured Memory: one sample at a time<sup>11</sup>







Figure: Sample-wise self-consistency and block-diagonal structures.



Figure: t-SNE of learned features . Left: U-CTRL and Right: MoCoV2.

EMP-SSL: Towards Self-Supervised Learning in One Training Epoch, Shengbang Tong, Yubei Chen, Yi Ma, Yann Lecun, arXiv:2304.03977

<sup>11</sup>Unsupervised Learning of Structured Representations via Closed-Loop Transcription, S. Tong, Yann LeCun, and Yi Ma, arXiv:2210.16782, 2022.

### **Towards AI 2.0: Time to Learn from Nature Again?**

Similar characteristics and mechanisms are ubiquitous in nature!

- Sparse coding in visual cortex (Olshausen, Nature 1996)<sup>12</sup>.
- Subspace embedding (Tsao, Cell 2017, Nature 2020).<sup>13</sup>
- Predictive coding in visual cortex (Rao, Nature Neuroscience 1999).

Article







- Face cells display flat tuning along dimensions orthogona the axis being coded
- The axis model is more efficient, robust, and flexible than t exemplar model
- Face patches ML/ME and AM carry complementar information about faces





#### **Towards AI 2.0 (Neural Science)**

#### A position paper about Intelligence in 2022:

### On the Principles of Parsimony and Self-Consistency for the Emergence of Intelligence

Yi Ma<sup>†‡1</sup>, Doris Tsao<sup>†2</sup>, Heung-Yeung Shum<sup>†3</sup>

<sup>1</sup>Electrical Engineering and Computer Science Department, University of California, Berkeley, CA 94720, USA <sup>2</sup>Department of Molecular & Cell Biology and Howard Hughes Medical Institute, University of California, Berkeley, CA 94720, USA <sup>3</sup>International Digital Economy Academy, Shenzhen 518045, China <sup>†</sup>E-mail: yima@eecs.berkeley.edu; dortsao@berkeley.edu; hshum@idea.edu.cn

Ma+TS, FITEE 2022

### **Towards AI 2.0 (Neural Science)**

- **Parsimony:** what's in neuroscience to verify this principle?
- Self-consistency: what's in neuroscience to verify this principle?
- Forward optimization versus backward propagation?
- Closed-loop versus open-loop?
- **Self-correcting** or self-improving mechanisms?

The Forward-Forward Algorithm for Training Deep Neural Networks

Invited Talk at NeurIPS 2022 Thurs 01 Dec 02:30 PM CST [ Hall H ]

PROCESSIN



## **Towards AI 2.0: How to Implement (Computer Science)**

To understand intelligence, one must understand computational complexity:



# Incomputable $\Rightarrow$ computable $\Rightarrow$ tractable $\Rightarrow$ scallable $\Rightarrow$ naturalKolmogorovTuring & ShannonNP vs PDNN & BPClosed-loop<br/>& feedback?

complexity

A massive parallel, distributed, hierarchical system of autoencoders (cortical columns in cortex)



### **Towards AI 2.0: Time to Learn from Nature Again?**



A unified purpose of intelligence: maximize "information gain" with every unit, at every stage!

### What is Intelligence?

**Definition [Intelligence]:** an intelligent system is one that has the mechanisms for **self-correcting** and **self-improving** its existing knowledge (or information).

$$\label{eq:Knowledge} \begin{split} & \text{Knowledge} = \int_0^t \text{Intelligence}, \\ & \text{Intelligence} = \frac{d}{dt} \text{Knowledge}. \end{split}$$

Any system without such mechanisms, however large, does not have any intelligence!





Who has intelligence, who has knowledge?

#### **Evolution of Intelligence in Nature: Four Stages**









Phylogenetic

Ontogenetic

Societal

**Artificial** 

#### Intelligence is all about how to encode and improve information for better prediction of the world!

	Phylogentic	Ontogenetic	Societal	Artificial
Codebook	Amino Acids	Neural Networks	Alphabet & Words   Mathematics/L	
Information	formation Genes/DNAs Memory		Natural Languages	Scientific Facts
Improvement	Natural Selection	Continuous Feedback	Trial & Error	Hypothesis Testing

### **Today's "Artificial Intelligence" is not that Artificial Intelligence!**

A quote from the 1956 Dartmouth proposal: "An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves".

	1940s (animal intelligence)	1956 (unique to human)	Today's Al (animal or human?)	
ERNETICS	<ul> <li>Signal processing</li> <li>Information Rep.</li> <li>Prediction</li> <li>Error correction</li> <li>Optimal control</li> <li>Game theory</li> </ul>	<ul> <li>Abstraction</li> <li>Logic deduction</li> <li>Causality</li> <li>Hypothesis forming &amp; testing</li> <li>Problem solving</li> </ul>	<ul> <li>Denoising</li> <li>Compression</li> <li>Object recognition</li> <li>Image generation</li> <li>Text generation</li> <li>Reinforce learning</li> </ul>	

IW.

#### An Open Question: Scientific Tests for Intelligence?

#### How to scientifically certify different types of ability of an "intelligent" system:

- Memorize: simply having memorized a large amount of knowledge-carrying data and regenerating them;
- Self-Learn: being able to autonomously and continuously develop better knowledge from new observations;
- Understand: having truly understood existing knowledge and knowing how to deduce and apply it correctly;
- **Theorize**: being able to generate new scientific hypotheses and mathematical theories and verify them.





#### Seek a scientific and theoretical foundation for Intelligence:

- what to learn? parsimony
- how to learn? compression
- why correct? consistency



"Everything should be made as simple as possible, but not any simpler."

-- Albert Einstein



#### Epilogue

### Coming soon: A new textbook in 2025

Learning Deep Representations of Data Distributions

Yi Ma, University of Hong Kong

Druv Pai, University of California, Berkeley

Sam Buchanan, Toyota Technological Institute at Chicago

Peng Wang, University of Macau

with contributions from:

Yaodong Yu, University of Maryland, College Park

#### **Acknowledgement**

#### A Truly Multi-University and Multi-Disciplinary Effort from Academia!







Sam Bu

Sam Buchanan TTI Chicago



Yaodong Yu OpenAl, U. Maryland



Ziyang Wu UC Berkeley



Shengbang Tong New York University



Tianzhe Chu Hong Kong University



Benjamin Haeffele Johns Hopkins



Peng Wang UMichigan



Simon Zhai UC Berkeley, DeepMind



Jack Bai

UIUC

Ryan Chan UPenn







#### HKU Musketeers Foundation Institute of Data Science 香港大學同心基金數據科學研究院



## **Pursuing the Nature of Intelligence**

## Thanks



