

# MI is All You Need

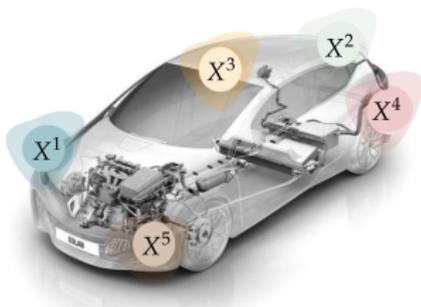
Understanding Complex Multivariate Systems Through the Lenses of GenAI

Prof. Pietro Michiardi

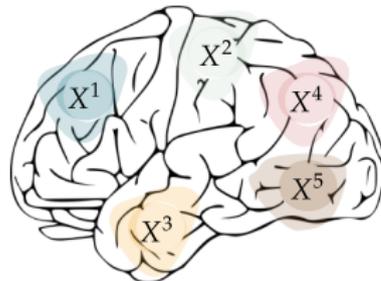


# Introduction

Complex systems are often described by **multivariate** information



Sensors



Brain regions

Understanding the **relationship** among multiple random variables is crucial to analyse **information content and flow** in these systems

# What do we use to study information?

- Shannon's Mutual Information (MI):  $\mathcal{I}(X^1; X^2)$
- Not interpretable for large systems  $X = \{X^1, \dots, X^N\}$ ,  $N > 3$

## PID

- Requires a partition into sources and one target
- Not scalable

## O-information

- No partition needed
- Scalable

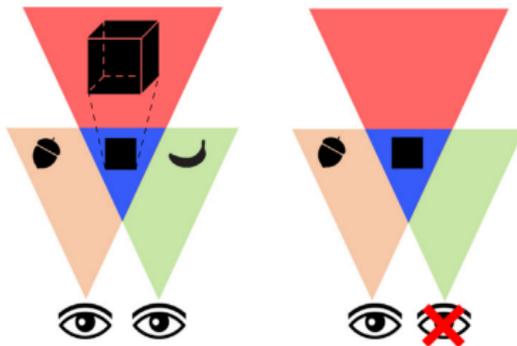
*SOTA is limited to discrete or Gaussian distributions*

Our methods, **MINDE** and **S $\Omega$ I**, estimate MI and O-information on **arbitrary continuous** systems of **any number** of variables

# Multivariate interactions

**Redundancy** : The **shared** information between variables, which can be recovered from variables or subset of variables

**Synergy** : The information that arises from **jointly** observing the variables but not accessible from individual variables alone



# O-information: a system-wide global measure

$$\boxed{\Omega(X) = \mathcal{T}(X) - \mathcal{D}(X)} \quad \left\{ \begin{array}{l} \Omega(X) > 0 \quad \textit{Redundancy} \\ \Omega(X) < 0 \quad \textit{Synergy} \end{array} \right.$$

$\mathcal{T}(X) \Rightarrow$  Information each variable  $X^i$  **shares** with others

$\mathcal{D}(X) \Rightarrow$  **Additional** information the variables  $X^i$  carry about part of the system, when the remaining part is known

Functions of **Mutual Information** (or Entropy)

Gradient of  $\Omega(X)$  captures **individual** flow of information

$$\boxed{\partial_i \Omega(X) = \Omega(X) - \Omega(X^{\setminus i})}$$

Hard for **high dimensional** and **continuous** distributions

# How to estimate MI? GenAI to the rescue!

Consider a **joint** generative diffusion model for two variables  $X^1$  and  $X^2$ :

$$\begin{cases} d[X_t^1, X_t^2]^\top = [X_t^1, X_t^2]^\top dt + \sqrt{2} [dW_t^1, dW_t^2]^\top, \\ [X_0^1, X_0^2] \sim p^{X^1, X^2} \end{cases}$$

In our work we show that the following holds:

$$\text{KL} [p^{X^1} \parallel p^{X^2}] = \mathbb{E}_{x \sim p^{X^1}} \left[ \int_0^T \underbrace{\left\| \nabla \log p_t^{X^1}(x) - \nabla \log p_t^{X^2}(x) \right\|^2}_{\text{Difference of score functions}} dt \right]$$

**Approximate** the KL by learning to denoise  $X_t^i$  :

$$\nabla \log p_t^{X^i}(x) = \frac{1}{2t} \left( \underbrace{\mathbb{E}[X^i | X_t^i]}_{\text{Denoiser} \approx_{\epsilon_\theta}(\cdot)} - x \right)$$

# MI estimation

So what? We have a way to estimate MI

$$\begin{aligned}\mathcal{I}(X^1; X^2) &=_{\text{KL}} [\rho(X^1, X^2) \parallel \rho(X^1)\rho(X^2)] \\ &= \mathbb{E}_{\rho(X^2)} [\text{KL} [\rho(X^1 | X^2) \parallel \rho(X^1)]]\end{aligned}$$

MINDE: **mutual information neural diffusion estimation** [ICLR 2024]<sup>1</sup>:

$$\mathcal{I}(X^1, X^2) = \int \frac{1}{4t^2} \mathbb{E} \|\mathbb{E}[X^1 | X_t^1] - \mathbb{E}[X^1 | X_t^1, X^2]\|^2 dt$$

Just compare the denoiser output when the variable  $X^1$  is denoised **alone** or **conditioned** on  $X^2$  !

---

<sup>1</sup>We have a discrete version of our MI estimator, preview at Delta Workshop, ICLR 2025

## $S\Omega$ : Score-based O-information estimation [ICML 2024]

Rewrite  $\mathcal{T}(X)$  and  $\mathcal{D}(X)$  in terms of KL divergence, apply previous results:

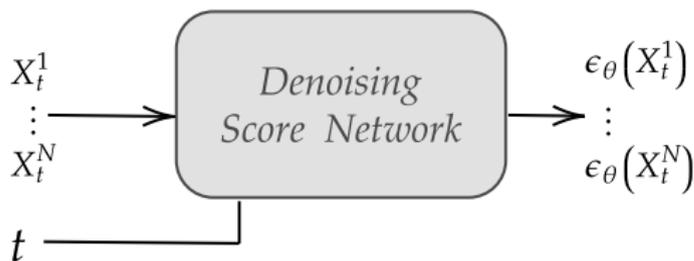
$$\mathcal{T}(X) = \int \frac{1}{4t^2} \mathbb{E} \left\| \mathbb{E}[X | \mathbf{X}_t] - [\mathbb{E}[X^i | \mathbf{X}_t^i]]_{i=1}^N \right\|^2 dt$$

Compare denoiser output when all the variables are denoised **jointly** or **marginally**

$$\mathcal{D}(X) = \int \frac{1}{4t^2} \mathbb{E} \left\| \mathbb{E}[X | \mathbf{X}_t] - [\mathbb{E}[X^i | \mathbf{X}_t^i, \mathbf{X}^{\setminus i}]]_{i=1}^N \right\|^2 dt$$

Compare the denoiser output when all the variables are denoised **jointly** or **conditionally** on the remaining clean variables

# Amortized approach using a unique network



---

## Algorithm 1: $\Sigma$ O-information estimation

---

**Input:**  $X = \{X^i\}_{i=1}^N, t \sim \mathcal{U}[0, T], X_t = X + \sqrt{2t}W$

$\epsilon(X_t) \leftarrow \epsilon_\theta([X_t^1, \dots, X_t^N], t)$  // Joint

**for**  $i = 1$  **to**  $N$  **do**

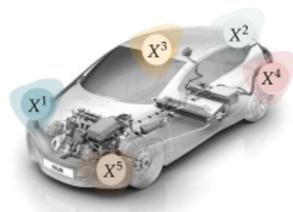
$\epsilon(X_t^i | X^{\setminus i}) \leftarrow \epsilon_\theta([X^1, \dots, X_t^i, \dots, X^N], t)$  // Conditional

$\epsilon(X_t^i) \leftarrow \epsilon_\theta(X_t^i, t)$  // Marginal

**Return**  $\underbrace{\frac{1}{4t^2} \left\| \epsilon(X_t) - [\epsilon(X_t^i)]_{i=1}^N \right\|^2}_{\mathcal{T}(X)} - \underbrace{\frac{1}{4t^2} \left\| \epsilon(X_t) - [\epsilon(X_t^i | X^{\setminus i})]_{i=1}^N \right\|^2}_{\mathcal{D}(X)}$

---

# Applications: Automotive (1)



- Sensors are **unreliable**: can we exploit redundancy?
- Objective: **cross generation**

How? MLD: multivariate generative model [ENTROPY 2024]

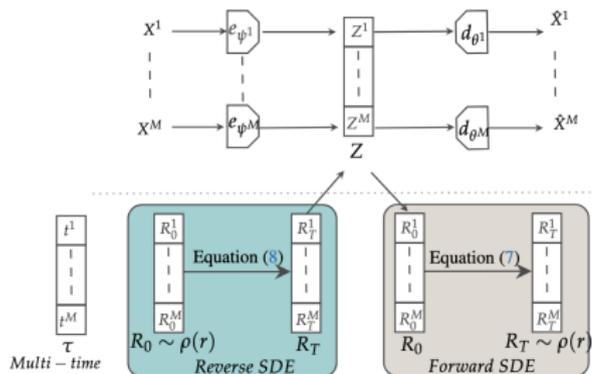
$$\begin{cases} d[X_t^1, X_t^2]^\top = [X_t^1, X_t^2]^\top dt + \sqrt{2} [dW_t^1, dW_t^2]^\top, \\ [X_0^1, X_0^2]^\top \sim p^{A,B}, \end{cases}$$

where  $X^1$  and  $X^2$  correspond to two **modalities** (for simplicity)

# Applications: Automotive (2)

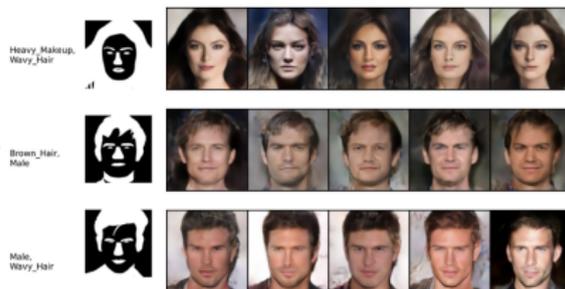
## Implementation

- Project inputs to latent variables
- Learn a joint diffusion model, use “multiple arrows of time”



## Results

- Dataset: 3 modalities (text, segmentation map, image)
- Any-to-any generation, **coherence** is defined by the joint score



# Applications: Alignment (1)

"A painting of an  
elephant with glasses"



- Common alignment issues: Catastrophic neglecting, Incorrect attribute binding, Incorrect spatial layout
- Existing solutions in the literature:
  - Test time: linguistic steering of generative pathways
  - Fine-tuning: ask GPT for help
- All methods require auxiliary LLM models

# Applications: Alignment (2)

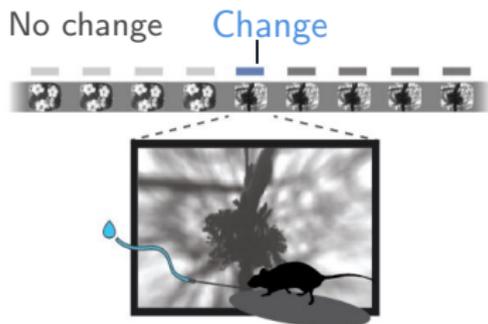
Our method: MITUNE [ICLR 2025]<sup>2</sup>

- Self-supervised fine-tuning: all is done with the generative model
- Steps:
  - Generate synthetic data using the pre-trained model
  - Compute **point-wise** MI for each prompt-image pair
  - Select top-k pairs with the highest MI
  - Fine-tune with adapters



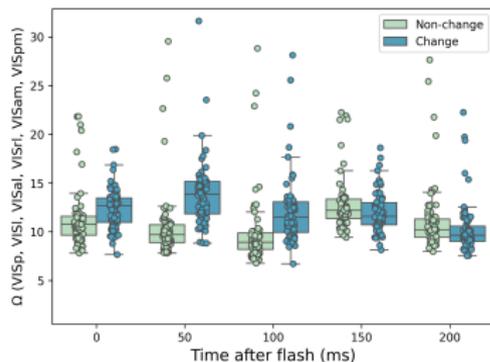
<sup>2</sup>We also have a version for Rectified Flows, preview at Delta Workshop, ICLR 2025.

# Applications: Neuroscience



O-information in the mice  
brain

$\Omega I$  is used to estimate  
O-information for each 50ms  
bin of spikes recording after the  
stimulus flash



6 brain regions

Higher **redundant** information in the  
visual cortex regions is transmitted in  
case of a flash with new scene

# Conclusion

- GenAI-based extension to Information Theory!
- New information measures unlock a wide array of applications
  - Neuroscience, Cellular development studies
  - Learning from unpaired data
  - Synthetic data augmentation, compression, ...
- Several fundamental problems related to neural estimation
  - Sample efficiency
  - Computational scalability
  - Interpretability and explainability

# Thank you !

Joint work with Prof. Giulio Franzese (EURECOM)

Published in ICLR 2024, ICML 2024, ENTROPY 2024, ICLR 2025

Supported by the Huawei Labs Paris, MUSE-COM2 CHIST-ERA project  
<https://musecom2.eu/>

# Backup Slides

# Why we need Mutual Information

Here's our complex, multivariate system, in abstract terms:

$$X = \underbrace{\{X^1, \dots, X^{i-1}\}}_{X^{<i}}, X^i, \underbrace{\{X^{i+1}, \dots, X^N\}}_{X^{>i}}, X^{\setminus i} = \{X^{<i}, X^{>i}\}$$

- **Total correlation:**  $\mathcal{T}(X) = \sum \mathcal{I}(X^i; X^{>i})$   
How much information each variable  $X^i$ , **shares** with  $X^{>i}$ , which suggests a *redundant* scenario
- **Dual total correlation:**  $\mathcal{D}(X) = \sum \mathcal{I}(X^i; X^{<i} | X^{>i})$   
How much **additional** information the variables  $X^i$  carry about  $X^{<i}$  if  $X^{>i}$  is also available which suggests a *synergistic* scenario

# Technical Details I

- Consider random variable  $A$  with probability measure  $p^A(x)dx$
- Build a simple SDE in  $[0, T]$  with initial conditions  $\sim p^A$

$$\begin{cases} dX_t = -X_t dt + \sqrt{2}dW_t, \\ X_0 \sim p^A \end{cases} \quad (1)$$

- This SDE corresponds to a path measure  $\mathbb{P}^A$
- It is possible to show that two SDEs which differ only by initial conditions have KL divergence

$$\text{KL}[\mathbb{P}^A \parallel \mathbb{P}^B] = \mathbb{E}_{\mathbb{P}^A} \left[ \log \frac{d\mathbb{P}^A}{d\mathbb{P}^B} \right] = \mathbb{E}_{\mathbb{P}^A} \left[ \log \frac{dp^A}{dp^B} \right] = \text{KL}[p^A \parallel p^B] \quad (2)$$

## Technical Details II

### Time-reversal $\hat{X}_t$ :

- KL between path measures is invariant to time-reversal  
 $\text{KL}[\mathbb{P}^A \parallel \mathbb{P}^B] = \text{KL}[\hat{\mathbb{P}}^A \parallel \hat{\mathbb{P}}^B]$
- Time reversal of SDE is again an SDE

$$d\hat{X}_t = \underbrace{\hat{X}_t + 2 \nabla \log p_{T-t}^A(\hat{X}_t)}_{\text{score function!}} dt + \sqrt{2} d\hat{W}_t \quad (3)$$

**Girsanov theorem:** (informal) express the KL between path measures corresponding to two SDEs with different drifts.

$$\text{KL}[\hat{\mathbb{P}}^{\mu^A} \parallel \hat{\mathbb{P}}^{\mu^B}] \simeq \mathbb{E}_{\mathbb{P}^{\mu^A}} \left[ \int_0^T \|\nabla \log p_t^A(X_t) - \nabla \log p_t^B(X_t)\|^2 dt \right] \quad (4)$$

Combining with  $\text{KL}[\mathbb{P}^A \parallel \mathbb{P}^B] = \text{KL}[\hat{\mathbb{P}}^A \parallel \hat{\mathbb{P}}^B]$  we can obtain a KL-estimator!

# Mutual Information Neural Diffusion Estimation

Mutual Information between two random variables  $A, B$  (many equivalent formulations):

$$I(A, B) = \text{KL} [p^{A,B} \parallel p^A p^B] \quad (5)$$

Idea: estimation using score functions! Two families diffusion processes: joint (J) and conditional (C)

$$\begin{cases} d[X_t, Y_t]^\top = -[X_t, Y_t]^\top dt + \sqrt{2} [dW_t, dW_t']^\top \\ [X_0, Y_0]^\top \sim p^{A,B} \end{cases}$$
$$\begin{cases} dX_t = -X_t dt + \sqrt{2} dW_t \\ X_0 \sim p^{A|B} \end{cases}$$