

# Contrastive Multiview Coding

---

Chengkun Li

[lichengkun0805@gmail.com](mailto:lichengkun0805@gmail.com)

## 1. Contrastive Learning

### 1.1. Concepts & Basic Idea

### 1.2. How to train?

## 2. CMC

### 2.1. Introduction

### 2.2. Concepts & Basic Idea of CMC

### 2.3. CMC with two views

### 2.4. CMC with Multiple Views

### 2.5. Why CMC?

### 2.6. Relationship with Mutual Information

## 3. Conclusions & Inspirations

## Contrastive Multiview Coding<sup>1</sup>

- Apply **Contrastive Learning** method
- Use **multiple** views from the same scene
- Aiming to learn the **representation** of data

---

<sup>1</sup>*Yonglong Tian et al. Contrastive Multiview Coding (MIT& Google)*

## Contrastive Multiview Coding<sup>1</sup>

- Apply **Contrastive Learning** method
- Use **multiple** views from the same scene
- Aiming to learn the **representation** of data

---

<sup>1</sup>*Yonglong Tian et al. Contrastive Multiview Coding (MIT& Google)*

## Contrastive Multiview Coding<sup>1</sup>

- Apply **Contrastive Learning** method
- Use **multiple** views from the same scene
- Aiming to learn the **representation** of data

---

<sup>1</sup>*Yonglong Tian et al. Contrastive Multiview Coding (MIT& Google)*

## Contrastive Multiview Coding<sup>1</sup>

- Apply **Contrastive Learning** method
- Use **multiple** views from the same scene
- Aiming to learn the **representation** of data

---

<sup>1</sup>*Yonglong Tian et al. Contrastive Multiview Coding (MIT& Google)*

## 1. Contrastive Learning

### 1.1. Concepts & Basic Idea

### 1.2. How to train?

## 2. CMC

### 2.1. Introduction

### 2.2. Concepts & Basic Idea of CMC

### 2.3. CMC with two views

### 2.4. CMC with Multiple Views

### 2.5. Why CMC?

### 2.6. Relationship with Mutual Information

## 3. Conclusions & Inspirations

## Definition of Contrastive Learning:

A learning paradigm that learns to tell *distinctiveness*



## Definition of Contrastive Learning:

A learning paradigm that learns to (by) tell (ing) *distinctiveness*

## Definition of Contrastive Learning:

A learning paradigm that learns to (by) tell (ing) *distinctiveness*

Puzzle...

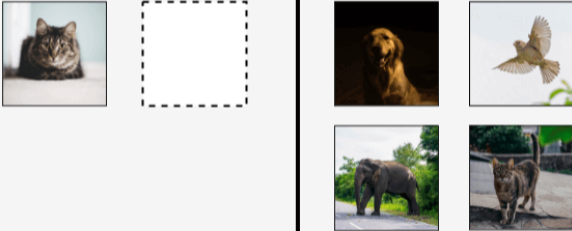
Match the correct animal



Can we teach machines this way?

Puzzle...

Match the correct animal



The puzzle consists of a 2x2 grid of images on the left and a 2x2 grid of images on the right, separated by a vertical line. The top-left image is a tabby cat. The top-right image is a dashed square. The bottom-left image is a golden retriever. The bottom-right image is a small white bird in flight. The bottom-left image is an elephant. The bottom-right image is a tabby cat.

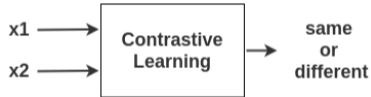
Can we teach machines this way?

# Problem Formulation

**Goal:** Teach machines to distinguish between **similar** and **dissimilar** things

What the machines need:

1. Similar & dissimilar data
2. Ability to represent the image (data)



3. Ability to quantify if two images are similar by their representation

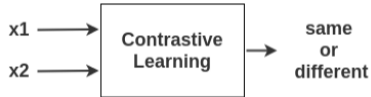
How?

# Problem Formulation

**Goal:** Teach machines to distinguish between **similar** and **dissimilar** things

**What the machines need:**

1. Similar & dissimilar data
2. Ability to *represent* the image (data)



3. Ability to *quantify* if two images are similar by their representation

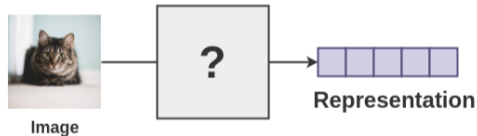
How?

# Problem Formulation

**Goal:** Teach machines to distinguish between **similar** and **dissimilar** things

**What the machines need:**

1. Similar & dissimilar data
2. Ability to *represent* the image (data)



3. Ability to *quantify* if two images are similar by their representation

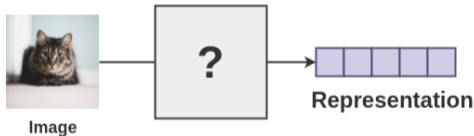
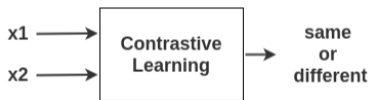
How?

# Problem Formulation

**Goal:** Teach machines to distinguish between **similar** and **dissimilar** things

**What the machines need:**

1. Similar & dissimilar data
2. Ability to *represent* the image (data)



3. Ability to *quantify* if two images are similar by their representation

How?

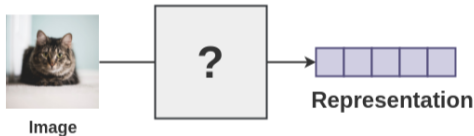
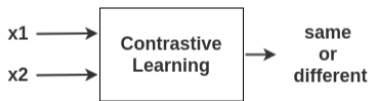


# Problem Formulation

**Goal:** Teach machines to distinguish between **similar** and **dissimilar** things

**What the machines need:**

1. Similar & dissimilar data
2. Ability to *represent* the image (data)

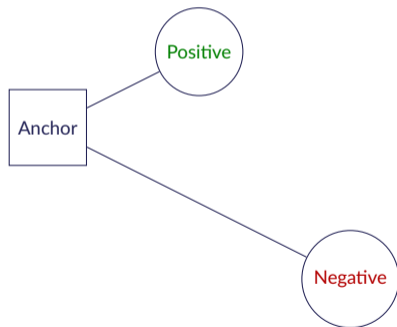


3. Ability to *quantify* if two images are similar by their representation

## How?

## Definitions

Anchor, Positive, Negative...



## Similar data

1. *Supervised way*: data with same label
2. *Unsupervised way/Self-Supervised way*:
  - data augmentation
  - Multi-modal
  - random cropping

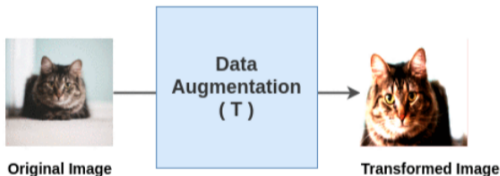
## Similar data

1. *Supervised way*: data with same label
2. *Unsupervised way/Self-Supervised way*:
  - data augmentation

- Multi modal
- random cropping

## Similar data

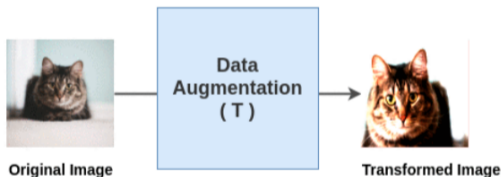
1. *Supervised way*: data with same label
2. *Unsupervised way/Self-Supervised way*:
  - data augmentation



- Multi modal
- random cropping

## Similar data

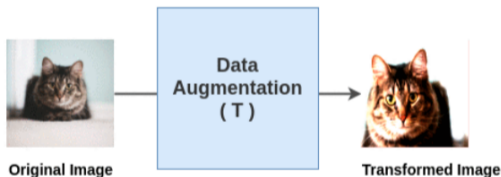
1. *Supervised way*: data with same label
2. *Unsupervised way/Self-Supervised way*:
  - data augmentation



- Multi modal
- random cropping

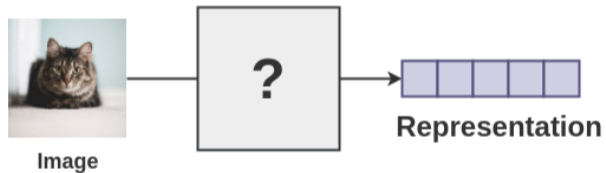
## Similar data

1. *Supervised way*: data with same label
2. *Unsupervised way/Self-Supervised way*:
  - data augmentation



- Multi modal
- random cropping

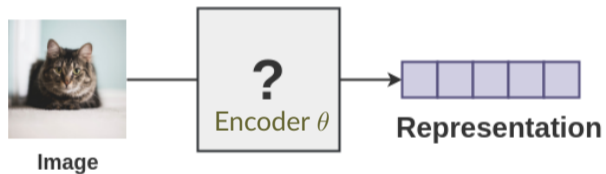
## Data Representation



$$z_i = f_{\theta}(x_i)$$

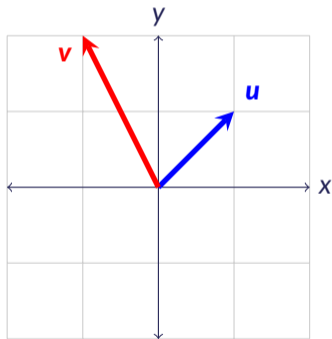


## Data Representation



$$z_i = f_{\theta}(x_i)$$

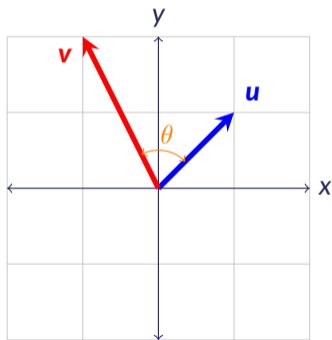
How to quantify the distinctiveness?



**Cosine Similarity**

$$\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|} \quad (1)$$

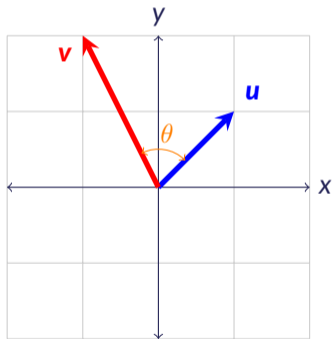
How to quantify the distinctiveness?



**Cosine Similarity**

$$\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|} \quad (1)$$

How to quantify the distinctiveness?



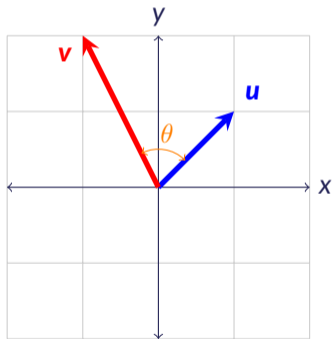
**Cosine Similarity**

$$\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|} \quad (1)$$

sometimes, add a dynamic range controlling hyperparam  $\tau$ :

$$\text{sim}(u, v) \cdot \frac{1}{\tau} = \frac{u^T v}{\|u\| \|v\| \tau}, \quad \tau \in [-1, 1] \quad (2)$$

How to quantify the distinctiveness?



**Cosine Similarity**

$$\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|} \quad (1)$$

sometimes, add a dynamic range controlling hyperparam  $\tau$ :

$$\text{sim}(u, v) \cdot \frac{1}{\tau} = \frac{u^T v}{\|u\| \|v\|^\tau}, \quad \tau \in [-1, 1] \quad (2)$$

( Some papers incorporate  $\tau$  into  $\text{sim}(u, v)$ )

## 1. Contrastive Learning

1.1. Concepts & Basic Idea

1.2. How to train?

## 2. CMC

2.1. Introduction

2.2. Concepts & Basic Idea of CMC

2.3. CMC with two views

2.4. CMC with Multiple Views

2.5. Why CMC?

2.6. Relationship with Mutual Information

## 3. Conclusions & Inspirations

# How to train?

## Steps (self-supervised)

1. Create **positive/negative** (in/congruent) data pairs (e.g. data augmentations)
2. Compute Data Representation (feature extraction)
3. Compute *contrastive loss* to optimize representation in previous step

# How to train?

## Steps (self-supervised)

1. Create **positive/negative** (in/congruent) data pairs (e.g. data augmentations)
2. Compute Data Representation (feature extraction)
3. Compute *contrastive loss* to optimize representation in previous step



# How to train?

## Steps (self-supervised)

1. Create **positive/negative** (in/congruent) data pairs (e.g. data augmentations)
2. Compute Data Representation (feature extraction)
3. Compute *contrastive loss* to optimize representation in previous step

# How to train?

## Steps (self-supervised)

1. Create **positive/negative** (in/congruent) data pairs (e.g. data augmentations)
2. Compute Data Representation (feature extraction)
3. Compute *contrastive loss* to optimize representation in previous step

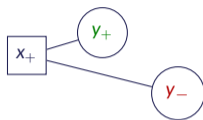
### Optimization intuition

Pull similar pairs **together**, push dissimilar pairs **away from** each others



# How to train?

## Common Loss functions



$x_+$ : Anchor;  $y_+$ : Positive;  $y_-$ : Negative

### 1. Triplet margin

$$\max \left( \|f(x_+) - f(y_+)\|^2 - \|f(x_+) - f(y_-)\|^2 + m, 0 \right) \quad (3)$$

### 2. NCE loss

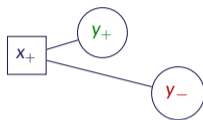
$$\log \sigma(\text{dis}(x_+, y_+) / \tau) + \log \sigma(-\text{dis}(x_+, y_-^i) / \tau) \quad (4)$$

### 3. k-pair loss (softmax-like)

$$-\log \frac{\exp(\text{sim}(x_+, y_+) / \tau)}{\exp(\text{sim}(x_+, y_+) / \tau) + \sum_{i=1}^k \exp(\text{sim}(x_+, y_-^i) / \tau)} \quad (5)$$

# How to train?

## Common Loss functions



$x_+$ : Anchor;  $y_+$ : Positive;  $y_-$ : Negative

### 1. Triplet margin

$$\max \left( \|f(x_+) - f(y_+)\|^2 - \|f(x_+) - f(y_-)\|^2 + m, 0 \right) \quad (3)$$

### 2. NCE loss

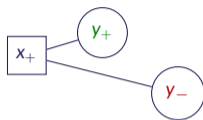
$$\log \sigma (\text{dis} (x_+, y_+) / \tau) + \log \sigma \left( -\text{dis} (x_+, y_-^i) / \tau \right) \quad (4)$$

### 3. k-pair loss (softmax-like)

$$-\log \frac{\exp (\text{sim} (x_+, y_+) / \tau)}{\exp (\text{sim} (x_+, y_+) / \tau) + \sum_{i=1}^k \exp (\text{sim} (x_+, y_-^i) / \tau)} \quad (5)$$

# How to train?

## Common Loss functions



$x_+$ : Anchor;  $y_+$ : Positive;  $y_-$ : Negative

### 1. Triplet margin

$$\max \left( \|f(x_+) - f(y_+)\|^2 - \|f(x_+) - f(y_-)\|^2 + m, 0 \right) \quad (3)$$

### 2. NCE loss

$$\log \sigma (\text{dis} (x_+, y_+) / \tau) + \log \sigma \left( - \text{dis} (x_+, y_-^i) / \tau \right) \quad (4)$$

### 3. k-pair loss (softmax-like)

$$- \log \frac{\exp (\text{sim} (x_+, y_+) / \tau)}{\exp (\text{sim} (x_+, y_+) / \tau) + \sum_{i=1}^k \exp (\text{sim} (x_+, y_-^i) / \tau)} \quad (5)$$

## 1. Contrastive Learning

### 1.1. Concepts & Basic Idea

### 1.2. How to train?

## 2. CMC

### 2.1. Introduction

### 2.2. Concepts & Basic Idea of CMC

### 2.3. CMC with two views

### 2.4. CMC with Multiple Views

### 2.5. Why CMC?

### 2.6. Relationship with Mutual Information

## 3. Conclusions & Inspirations

## 1. Contrastive Learning

### 1.1. Concepts & Basic Idea

### 1.2. How to train?

## 2. CMC

### 2.1. Introduction

### 2.2. Concepts & Basic Idea of CMC

### 2.3. CMC with two views

### 2.4. CMC with Multiple Views

### 2.5. Why CMC?

### 2.6. Relationship with Mutual Information

## 3. Conclusions & Inspirations

## Contrastive Multiview Coding<sup>2</sup>

- Apply Contrastive Learning method
- Use **multiple** views from the same scene
- Aiming to learn the **representation** of data

---

<sup>2</sup>*Yonglong Tian et al. Contrastive Multiview Coding*



## Contrastive Multiview Coding<sup>2</sup>

- Apply Contrastive Learning method
- Use **multiple** views from the same scene
- Aiming to learn the **representation** of data

---

<sup>2</sup>*Yonglong Tian et al. Contrastive Multiview Coding*

## Contrastive Multiview Coding<sup>2</sup>

- Apply Contrastive Learning method
- Use **multiple** views from the same scene
- Aiming to learn the **representation** of data

---

<sup>2</sup>*Yonglong Tian et al. Contrastive Multiview Coding*

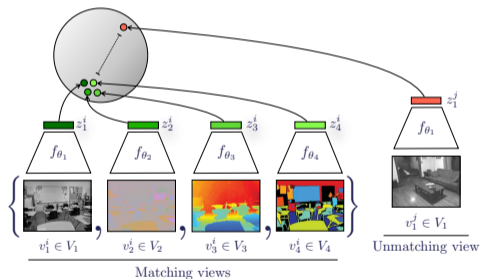
## Contrastive Multiview Coding<sup>2</sup>

- Apply Contrastive Learning method
- Use **multiple** views from the same scene
- Aiming to learn the **representation** of data

---

<sup>2</sup>*Yonglong Tian et al. Contrastive Multiview Coding*

Comparing with the previous self-supervised example,



CMC...

1. data augmentation  $\rightarrow$  views
2. two sample data  $\rightarrow$  multiple views

## 1. Contrastive Learning

### 1.1. Concepts & Basic Idea

### 1.2. How to train?

## 2. CMC

### 2.1. Introduction

### 2.2. Concepts & Basic Idea of CMC

### 2.3. CMC with two views

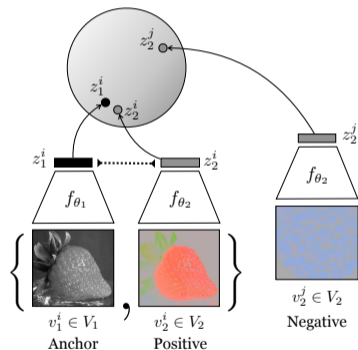
### 2.4. CMC with Multiple Views

### 2.5. Why CMC?

### 2.6. Relationship with Mutual Information

## 3. Conclusions & Inspirations

## Notations



### Recall: sim function

$$\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|}$$

- **Views:**  $V_1, V_2, \dots, V_N$

- **Data in a view:**  $v_a^b$

- **Positives:**  $\{v_1^i, v_2^i\}_{i=1}^N$

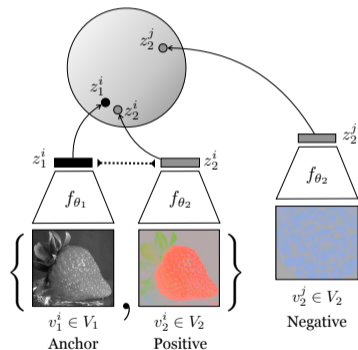
- **Negatives:**  $\{v_1^i, v_2^j\}, i, j \in N$

- **Discriminating function:**  $h_{\theta}(\cdot)$

$$h_{\theta}(\{v_1, v_2\}) = \exp\left(\frac{f_{\theta_1}(v_1) \cdot f_{\theta_2}(v_2)}{\|f_{\theta_1}(v_1)\| \cdot \|f_{\theta_2}(v_2)\|} \cdot \frac{1}{\tau}\right) \quad (6)$$

- **Encoders:**  $\theta_1, \theta_2, \dots, \theta_N$

## Notations



## Recall: sim function

$$\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|}$$

- **Views:**  $V_1, V_2, \dots, V_N$

- **Data in a view:**  $v_a^b$

- **Positives:**  $\{v_1^i, v_2^i\}_{i=1}^N$

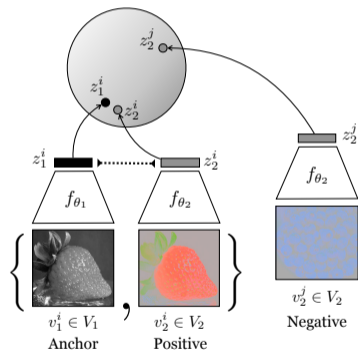
- **Negatives:**  $\{v_1^i, v_2^j\}, i, j \in N$

- **Discriminating function:**  $h_{\theta}(\cdot)$

$$h_{\theta}(\{v_1, v_2\}) = \exp\left(\frac{f_{\theta_1}(v_1) \cdot f_{\theta_2}(v_2)}{\|f_{\theta_1}(v_1)\| \cdot \|f_{\theta_2}(v_2)\|} \cdot \frac{1}{\tau}\right) \quad (6)$$

- **Encoders:**  $\theta_1, \theta_2, \dots, \theta_N$

## Notations



## Recall: sim function

$$\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|}$$

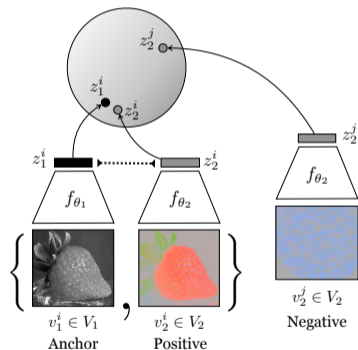
- **Views:**  $V_1, V_2, \dots, V_N$
- **Data in a view:**  $v_a^b$
- **Positives:**  $\{v_1^i, v_2^i\}_{i=1}^N$
- **Negatives:**  $\{v_1^i, v_2^j\}, i, j \in N$
- **Discriminating function:**  $h_{\theta}(\cdot)$

$$h_{\theta}(\{v_1, v_2\}) = \exp\left(\frac{f_{\theta_1}(v_1) \cdot f_{\theta_2}(v_2)}{\|f_{\theta_1}(v_1)\| \cdot \|f_{\theta_2}(v_2)\|} \cdot \frac{1}{\tau}\right) \quad (6)$$

- **Encoders:**  $\theta_1, \theta_2, \dots, \theta_N$



## Notations



## Recall: sim function

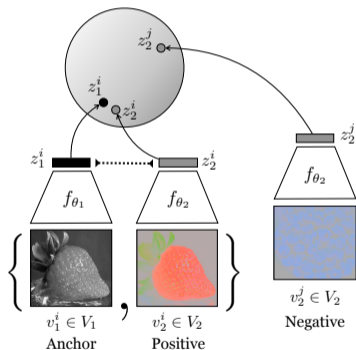
$$\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|}$$

- **Views:**  $V_1, V_2, \dots, V_N$
- **Data in a view:**  $v_a^b$
- **Positives:**  $\{v_1^i, v_2^i\}_{i=1}^N$
- **Negatives:**  $\{v_1^i, v_2^j\}, i, j \in N$
- **Discriminating function:**  $h_{\theta}(\cdot)$

$$h_{\theta}(\{v_1, v_2\}) = \exp\left(\frac{f_{\theta_1}(v_1) \cdot f_{\theta_2}(v_2)}{\|f_{\theta_1}(v_1)\| \cdot \|f_{\theta_2}(v_2)\|} \cdot \frac{1}{\tau}\right) \quad (6)$$

- **Encoders:**  $\theta_1, \theta_2, \dots, \theta_N$

## Notations



### Recall: sim function

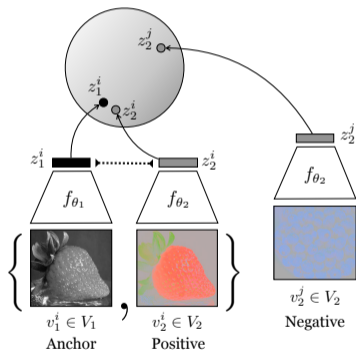
$$\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|}$$

- **Views:**  $V_1, V_2, \dots, V_N$
- **Data in a view:**  $v_a^b$
- **Positives:**  $\{v_1^i, v_2^i\}_{i=1}^N$
- **Negatives:**  $\{v_1^i, v_2^j\}, i, j \in N$
- **Discriminating function:**  $h_{\theta}(\cdot)$

$$h_{\theta}(\{v_1, v_2\}) = \exp\left(\frac{f_{\theta_1}(v_1) \cdot f_{\theta_2}(v_2)}{\|f_{\theta_1}(v_1)\| \cdot \|f_{\theta_2}(v_2)\|} \cdot \frac{1}{\tau}\right) \quad (6)$$

- **Encoders:**  $\theta_1, \theta_2, \dots, \theta_N$

## Notations



### Recall: sim function

$$\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|}$$

- **Views:**  $V_1, V_2, \dots, V_N$
- **Data in a view:**  $v_a^b$
- **Positives:**  $\{v_1^i, v_2^i\}_{i=1}^N$
- **Negatives:**  $\{v_1^i, v_2^j\}, i, j \in N$
- **Discriminating function:**  $h_{\theta}(\cdot)$

$$h_{\theta}(\{v_1, v_2\}) = \exp\left(\frac{f_{\theta_1}(v_1) \cdot f_{\theta_2}(v_2)}{\|f_{\theta_1}(v_1)\| \cdot \|f_{\theta_2}(v_2)\|} \cdot \frac{1}{\tau}\right) \quad (6)$$

- **Encoders:**  $\theta_1, \theta_2, \dots, \theta_N$

## 1. Contrastive Learning

### 1.1. Concepts & Basic Idea

### 1.2. How to train?

## 2. CMC

### 2.1. Introduction

### 2.2. Concepts & Basic Idea of CMC

### 2.3. CMC with two views

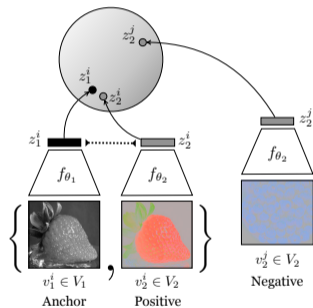
### 2.4. CMC with Multiple Views

### 2.5. Why CMC?

### 2.6. Relationship with Mutual Information

## 3. Conclusions & Inspirations

# CMC with two views



## k-pair loss

$$\ell = -\log \frac{\exp(\text{sim}(x_+, y_+) / \tau)}{\exp(\text{sim}(x_+, y_+) / \tau) + \sum_{i=1}^k \exp(\text{sim}(x_+, y_-^i) / \tau)}$$

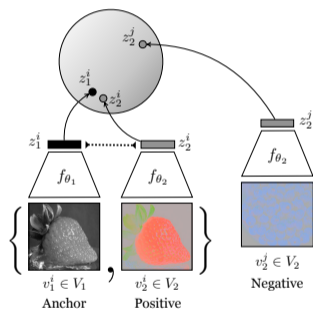
Contrastive loss between  $V_1$  (anchor view),  $V_2$

$$\mathcal{L}_{\text{contrast}}^{V_1, V_2} = -\mathbb{E}_{\{v_1^1, v_2^1, \dots, v_2^{k+1}\}} \left[ \log \frac{h_{\theta}(\{v_1^1, v_2^1\})}{\sum_{j=1}^{k+1} h_{\theta}(\{v_1^1, v_2^j\})} \right] \quad (7)$$

Total loss for two views:

$$\mathcal{L}(V_1, V_2) = \mathcal{L}_{\text{contrast}}^{V_1, V_2} + \mathcal{L}_{\text{contrast}}^{V_2, V_1} \quad (8)$$

# CMC with two views



## k-pair loss

$$\ell = -\log \frac{\exp(\text{sim}(x_+, y_+) / \tau)}{\exp(\text{sim}(x_+, y_+) / \tau) + \sum_{i=1}^k \exp(\text{sim}(x_+, y_-^i) / \tau)}$$

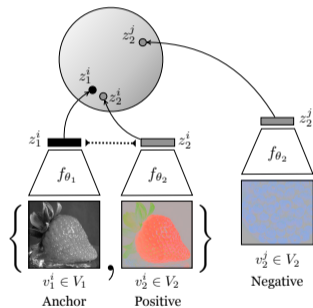
Contrastive loss between  $V_1$  (anchor view),  $V_2$

$$\mathcal{L}_{\text{contrast}}^{V_1, V_2} = - \mathbb{E}_{\{v_1^1, v_2^1, \dots, v_2^{k+1}\}} \left[ \log \frac{h_{\theta}(\{v_1^1, v_2^1\})}{\sum_{j=1}^{k+1} h_{\theta}(\{v_1^1, v_2^j\})} \right] \quad (7)$$

Total loss for two views:

$$\mathcal{L}(V_1, V_2) = \mathcal{L}_{\text{contrast}}^{V_1, V_2} + \mathcal{L}_{\text{contrast}}^{V_2, V_1} \quad (8)$$

# CMC with two views



## k-pair loss

$$\ell = -\log \frac{\exp(\text{sim}(x_+, y_+) / \tau)}{\exp(\text{sim}(x_+, y_+) / \tau) + \sum_{i=1}^k \exp(\text{sim}(x_+, y_-^i) / \tau)}$$

Contrastive loss between  $V_1$  (anchor view),  $V_2$

$$\mathcal{L}_{\text{contrast}}^{V_1, V_2} = - \mathbb{E}_{\{v_1^1, v_2^1, \dots, v_2^{k+1}\}} \left[ \log \frac{h_{\theta}(\{v_1^1, v_2^1\})}{\sum_{j=1}^{k+1} h_{\theta}(\{v_1^1, v_2^j\})} \right] \quad (7)$$

Total loss for two views:

$$\mathcal{L}(V_1, V_2) = \mathcal{L}_{\text{contrast}}^{V_1, V_2} + \mathcal{L}_{\text{contrast}}^{V_2, V_1} \quad (8)$$

## 1. Contrastive Learning

### 1.1. Concepts & Basic Idea

### 1.2. How to train?

## 2. CMC

### 2.1. Introduction

### 2.2. Concepts & Basic Idea of CMC

### 2.3. CMC with two views

### **2.4. CMC with Multiple Views**

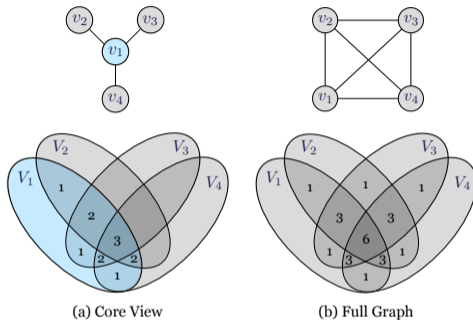
### 2.5. Why CMC?

### 2.6. Relationship with Mutual Information

## 3. Conclusions & Inspirations



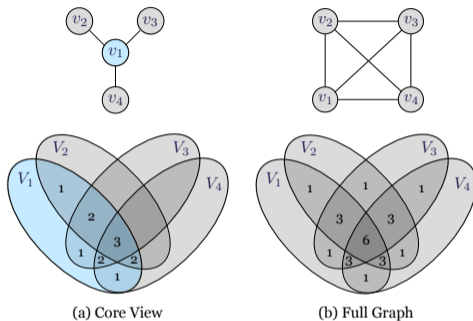
# CMC with Multiple Views



$$\mathcal{L}_C = \sum_{j=2}^M \mathcal{L}(V_1, V_j) \quad (9)$$

$$\mathcal{L}_F = \sum_{1 \leq i < j \leq M} \mathcal{L}(V_i, V_j) \quad (10)$$

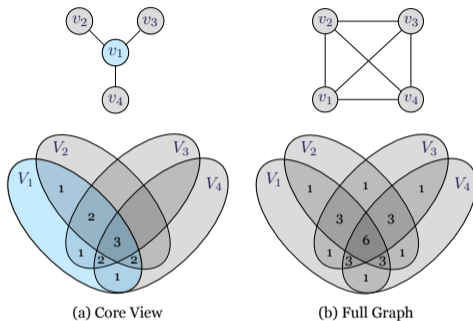
# CMC with Multiple Views



$$\mathcal{L}_C = \sum_{j=2}^M \mathcal{L}(V_1, V_j) \quad (9)$$

$$\mathcal{L}_F = \sum_{1 \leq i < j \leq M} \mathcal{L}(V_i, V_j) \quad (10)$$

# CMC with Multiple Views



$$\mathcal{L}_C = \sum_{j=2}^M \mathcal{L}(V_1, V_j) \quad (9)$$

$$\mathcal{L}_F = \sum_{1 \leq i < j \leq M} \mathcal{L}(V_i, V_j) \quad (10)$$

## 1. Contrastive Learning

### 1.1. Concepts & Basic Idea

### 1.2. How to train?

## 2. CMC

### 2.1. Introduction

### 2.2. Concepts & Basic Idea of CMC

### 2.3. CMC with two views

### 2.4. CMC with Multiple Views

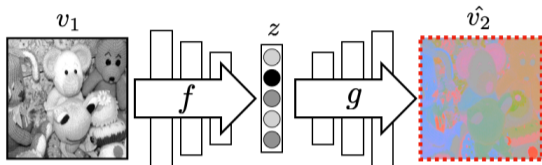
### 2.5. Why CMC?

### 2.6. Relationship with Mutual Information

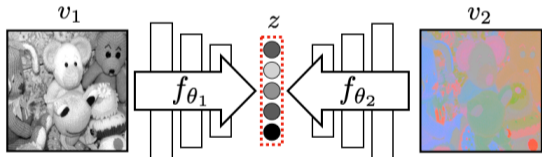
## 3. Conclusions & Inspirations

- CMC vs Encoder-Decoder
- CMC vs Supervised classifier

# CMC vs Encoder-Decoder



(a) Predictive learning



(b) Contrastive learning

For **predictive** learning:

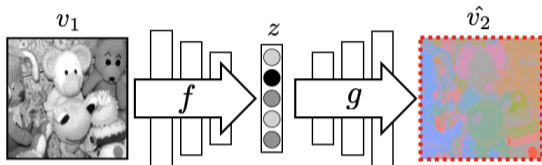
1.  $z = f(v_1)$  and  $\hat{v}_2 = g(z)$
2. Train  $f, g$  to make  $\hat{v}_2$  **closer to**  $v_2$  (at output space)

For **contrastive** learning:

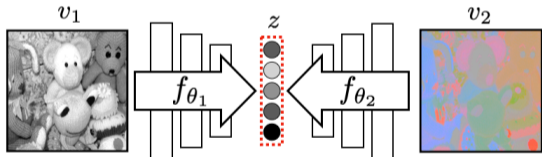
1.  $z_1 = f_{\theta_1}(v_1), z_2 = f_{\theta_2}(v_2)$
2. Train  $f_{\theta_1}, f_{\theta_2}$  to make  $z_1$  **closer to**  $z_2$  (at latent space)

Pixel wise  $\mathcal{L}_1, \mathcal{L}_2$  loss model complex structure poorly

# CMC vs Encoder-Decoder



(a) Predictive learning



(b) Contrastive learning

For **predictive** learning:

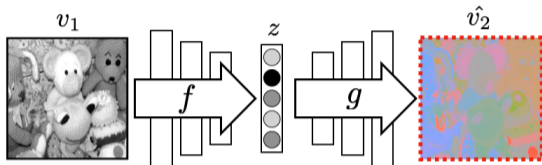
1.  $z = f(v_1)$  and  $\hat{v}_2 = g(z)$
2. Train  $f, g$  to make  $\hat{v}_2$  **closer to**  $v_2$  (at output space)

For **contrastive** learning:

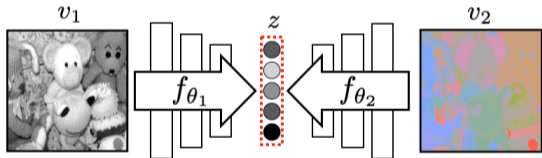
1.  $z_1 = f_{\theta_1}(v_1), z_2 = f_{\theta_2}(v_2)$
2. Train  $f_{\theta_1}, f_{\theta_2}$  to make  $z_1$  **closer to**  $z_2$  (at latent space)

Pixel wise  $\mathcal{L}_1, \mathcal{L}_2$  loss model complex structure poorly

# CMC vs Encoder-Decoder



(a) Predictive learning



(b) Contrastive learning

For **predictive** learning:

1.  $z = f(v_1)$  and  $\hat{v}_2 = g(z)$
2. Train  $f, g$  to make  $\hat{v}_2$  **closer to**  $v_2$  (at output space)

For **contrastive** learning:

1.  $z_1 = f_{\theta_1}(v_1), z_2 = f_{\theta_2}(v_2)$
2. Train  $f_{\theta_1}, f_{\theta_2}$  to make  $z_1$  **closer to**  $z_2$  (at latent space)

Pixel wise  $\mathcal{L}_1, \mathcal{L}_2$  loss model complex structure poorly



## Example



## Senses

- **Vision:** dog with a mask?
- **Acoustic:** "bark, bark"
- **Texture:** furry

Mutual Information

## Example



## Senses

- **Vision:** dog with a mask?
- **Acoustic:** "bark, bark"
- **Texture:** furry

Mutual Information

## Example



## Senses

- **Vision:** dog with a mask?
- **Acoustic:** "bark, bark"
- **Texture:** furry

Mutual Information

## Example



## Senses

- **Vision:** dog with a mask?
- **Acoustic:** "bark, bark"
- **Texture:** furry

Mutual Information

## Example



## Senses

- **Vision:** dog with a mask?
- **Acoustic:** "bark, bark"
- **Texture:** furry

## Mutual Information

## 1. Contrastive Learning

### 1.1. Concepts & Basic Idea

### 1.2. How to train?

## 2. CMC

### 2.1. Introduction

### 2.2. Concepts & Basic Idea of CMC

### 2.3. CMC with two views

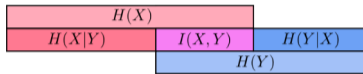
### 2.4. CMC with Multiple Views

### 2.5. Why CMC?

### 2.6. Relationship with Mutual Information

## 3. Conclusions & Inspirations

## Mutual Information<sup>3</sup>

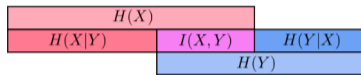


$$I(X; Y) = H(X) - H(X|Y) \quad (11)$$

$$I(X; Y) = \int_y \int_x p_{(X,Y)}(x, y) \log \left( \frac{p_{(X,Y)}(x, y)}{p_X(x)p_Y(y)} \right) dx dy \quad (12)$$

<sup>3</sup>Visual Information Theory *Christopher Olah*

## Mutual Information<sup>3</sup>



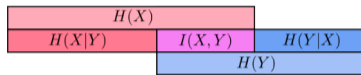
$$I(X; Y) = H(X) - H(X|Y) \quad (11)$$

$$I(X; Y) = \int_y \int_x p_{(X,Y)}(x, y) \log \left( \frac{p_{(X,Y)}(x, y)}{p_X(x)p_Y(y)} \right) dx dy \quad (12)$$

<sup>3</sup>Visual Information Theory *Christopher Olah*



## Mutual Information<sup>3</sup>



$$I(X; Y) = H(X) - H(X|Y) \quad (11)$$

$$I(X; Y) = \int_{\mathcal{Y}} \int_{\mathcal{X}} p_{(X,Y)}(x, y) \log \left( \frac{p_{(X,Y)}(x, y)}{p_X(x)p_Y(y)} \right) dx dy \quad (12)$$

<sup>3</sup>Visual Information Theory Christopher Olah

# Relationship with Mutual Information

## Mutual Information

$$I(X; Y) = \int_{\mathcal{Y}} \int_{\mathcal{X}} p_{(X,Y)}(x, y) \log \left( \frac{p_{(X,Y)}(x, y)}{p_X(x)p_Y(y)} \right) dx dy$$

Relationship between  $\mathcal{L}_{contrast}$  and MI,  $k$  is number of negatives

$$\begin{aligned} \mathcal{L}_{contrast} &\geq \log(k) - \mathbb{E}_{(z_1, z_2) \sim p_{z_1, z_2}(\cdot)} \log \left[ \frac{p(z_1, z_2)}{p(z_1)p(z_2)} \right] \\ &= \log(k) - I(z_1; z_2) \end{aligned}$$

we get<sup>4</sup>

$$I(v_i; v_j) \geq I(z_i; z_j) \geq \log(k) - \mathcal{L}_{contrast}$$

---

<sup>4</sup>Similar idea also appears in *Oord et al. Representation learning with contrastive predictive coding*

# Relationship with Mutual Information

## Mutual Information

$$I(X; Y) = \int_{\mathcal{Y}} \int_{\mathcal{X}} p_{(X,Y)}(x, y) \log \left( \frac{p_{(X,Y)}(x, y)}{p_X(x)p_Y(y)} \right) dx dy$$

Relationship between  $\mathcal{L}_{contrast}$  and MI,  $k$  is number of negatives

$$\begin{aligned} \mathcal{L}_{contrast} &\geq \log(k) - \mathbb{E}_{(z_1, z_2) \sim p_{z_1, z_2}(\cdot)} \log \left[ \frac{p(z_1, z_2)}{p(z_1)p(z_2)} \right] \\ &= \log(k) - I(z_1; z_2) \end{aligned}$$

we get<sup>4</sup>

$$I(v_i; v_j) \geq I(z_i; z_j) \geq \log(k) - \mathcal{L}_{contrast}$$

<sup>4</sup>Similar idea also appears in *Oord et al. Representation learning with contrastive predictive coding*

# Relationship with Mutual Information

## Mutual Information

$$I(X; Y) = \int_{\mathcal{Y}} \int_{\mathcal{X}} p_{(X,Y)}(x, y) \log \left( \frac{p_{(X,Y)}(x, y)}{p_X(x)p_Y(y)} \right) dx dy$$

Relationship between  $\mathcal{L}_{contrast}$  and MI,  $k$  is number of negatives

$$\begin{aligned} \mathcal{L}_{contrast} &\geq \log(k) - \mathbb{E}_{(z_1, z_2) \sim p_{z_1, z_2}(\cdot)} \log \left[ \frac{p(z_1, z_2)}{p(z_1)p(z_2)} \right] \\ &= \log(k) - I(z_1; z_2) \end{aligned}$$

we get<sup>4</sup>

$$I(v_i; v_j) \geq I(z_i; z_j) \geq \log(k) - \mathcal{L}_{contrast}$$

<sup>4</sup>Similar idea also appears in Oord et al. Representation learning with contrastive predictive coding

Better performance with higher MI?

$$I(z_i; z_j) \geq \log(k) - \mathcal{L}_{\text{contrast}}$$

Better performance with higher MI?

$$I(z_i; z_j) \geq \log(k) - \mathcal{L}_{\text{contrast}}$$

Better performance with higher MI?

$$I(z_i; z_j) \geq \log(k) - \mathcal{L}_{\text{contrast}}$$

**Not exactly**

# Reducing $I(v_1; v_2)$ with Spatial Distance



## Experiment setup:

- Two patches start at  $(x, y)$  and  $(x + d, y + d)$
- Train a linear classifier on pre-trained CMC representations
- Test classification accuracy



# Reducing $I(v_1; v_2)$ with Spatial Distance



## Experiment setup:

- Two patches start at  $(x, y)$  and  $(x + d, y + d)$
- Train a linear classifier on pre-trained CMC representations
- Test classification accuracy

# Reducing $I(v_1; v_2)$ with Spatial Distance



## Experiment setup:

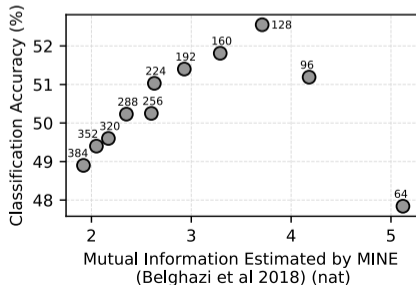
- Two patches start at  $(x, y)$  and  $(x + d, y + d)$
- Train a linear classifier on pre-trained CMC representations
- Test classification accuracy

# Reducing $I(v_1; v_2)$ with Spatial Distance



## Experiment setup:

- Two patches start at  $(x, y)$  and  $(x + d, y + d)$
- Train a linear classifier on pre-trained CMC representations
- Test classification accuracy



# Relationship with Mutual Information

Why?<sup>5</sup>

downstream task  $y$ , data  $x$

Too much information introduces **task-irrelevant** noise (nuisance)

---

<sup>1</sup>*Yonglong Tian et al. What Makes for Good Views for Contrastive Learning?*

# Relationship with Mutual Information

Why?<sup>5</sup>

downstream task  $y$ , data  $x$

Too much information introduces **task-irrelevant** noise (nuisance)

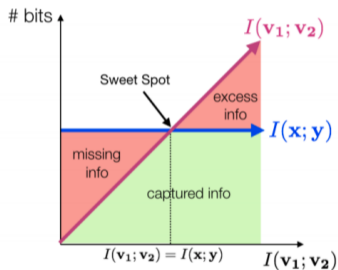
---

<sup>1</sup>*Yonglong Tian et al. What Makes for Good Views for Contrastive Learning?*

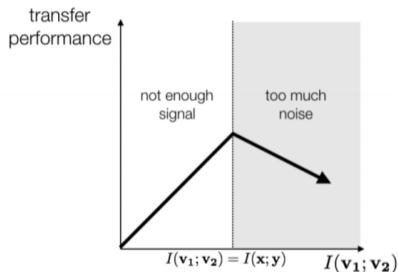
# Relationship with Mutual Information

Why?<sup>5</sup>

downstream task  $y$ , data  $x$



hypothesis



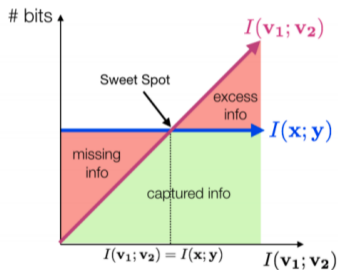
Too much information introduces **task-irrelevant** noise (nuisance)

<sup>1</sup>Yonglong Tian et al. What Makes for Good Views for Contrastive Learning?

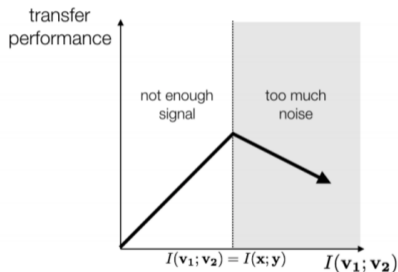
# Relationship with Mutual Information

Why?<sup>5</sup>

downstream task  $y$ , data  $x$



hypothesis



Too much information introduces **task-irrelevant** noise (nuisance)

<sup>1</sup>Yonglong Tian et al. What Makes for Good Views for Contrastive Learning?

## 1. Contrastive Learning

1.1. Concepts & Basic Idea

1.2. How to train?

## 2. CMC

2.1. Introduction

2.2. Concepts & Basic Idea of CMC

2.3. CMC with two views

2.4. CMC with Multiple Views

2.5. Why CMC?

2.6. Relationship with Mutual Information

## 3. Conclusions & Inspirations



# Conclusions & Inspirations

- CMC enables the learning of unsupervised representations from **multiple views** of datasets
- Mutual Information between views/modalities should be considered
- Contrastive Methods could be leveraged with/instead of Predictive Methods<sup>6</sup>
- ...

---

<sup>6</sup>Contrastive Learning for Unpaired Image-to-Image Translation (ECCV 2020)

# Conclusions & Inspirations

- CMC enables the learning of unsupervised representations from **multiple views** of datasets
- Mutual Information between views/modalities should be considered
- Contrastive Methods could be leveraged with/instead of Predictive Methods<sup>6</sup>
- ...

---

<sup>6</sup>Contrastive Learning for Unpaired Image-to-Image Translation (ECCV 2020)

# Conclusions & Inspirations

- CMC enables the learning of unsupervised representations from **multiple views** of datasets
- Mutual Information between views/modalities should be considered
- Contrastive Methods could be leveraged with/instead of Predictive Methods<sup>6</sup>
- ...

---

<sup>6</sup>Contrastive Learning for Unpaired Image-to-Image Translation (ECCV 2020)

# Conclusions & Inspirations

- CMC enables the learning of unsupervised representations from **multiple views** of datasets
- Mutual Information between views/modalities should be considered
- Contrastive Methods could be leveraged with/instead of Predictive Methods<sup>6</sup>
- ...

---

<sup>6</sup>Contrastive Learning for Unpaired Image-to-Image Translation (ECCV 2020)

**Thank you!**