# CLIP embedding analysis

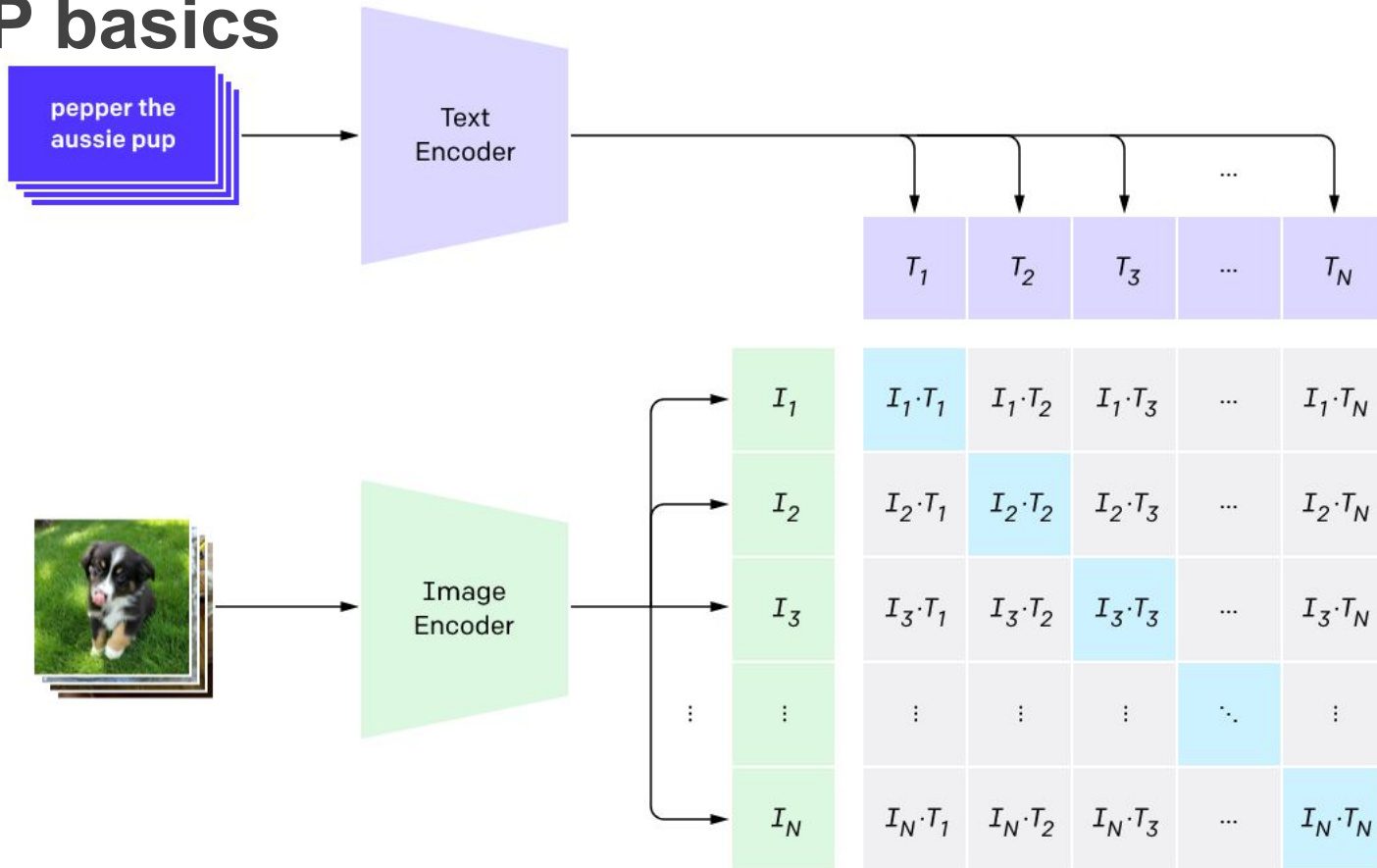**Semantic extraction and synthesis**

# Motivation

- Multi-modality approaches are of big interest of the SOTA research.
- Dimensionality reduction is important task in ML/DL.
- Managed synthesis is key issue in generative models applications.

# **Tasks**

- dimensionality reduction
- embedding clustering
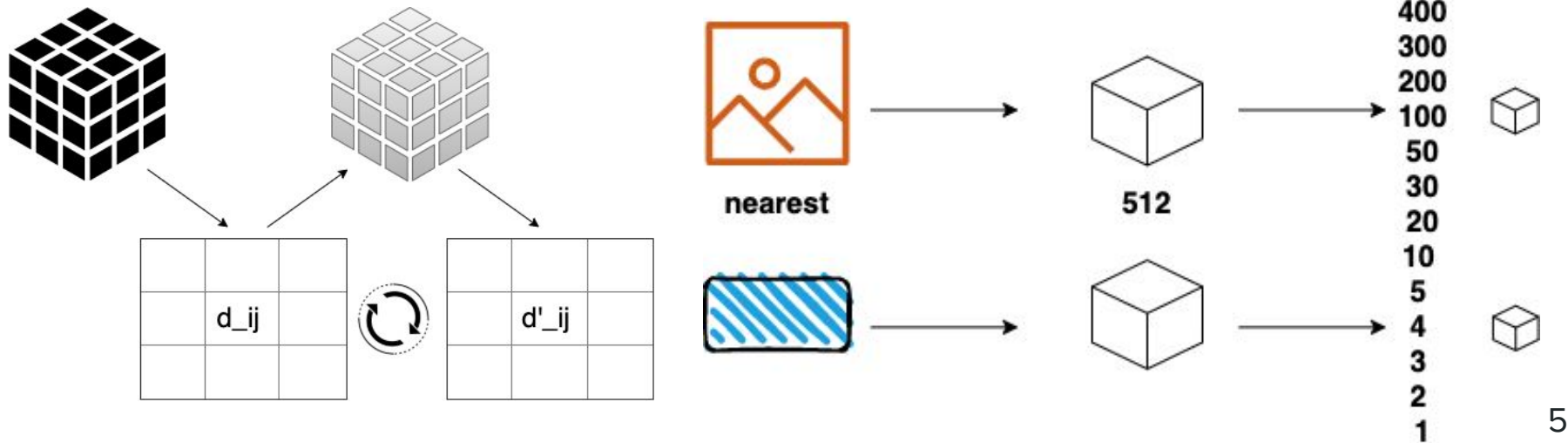- disentanglement and managed synthesis
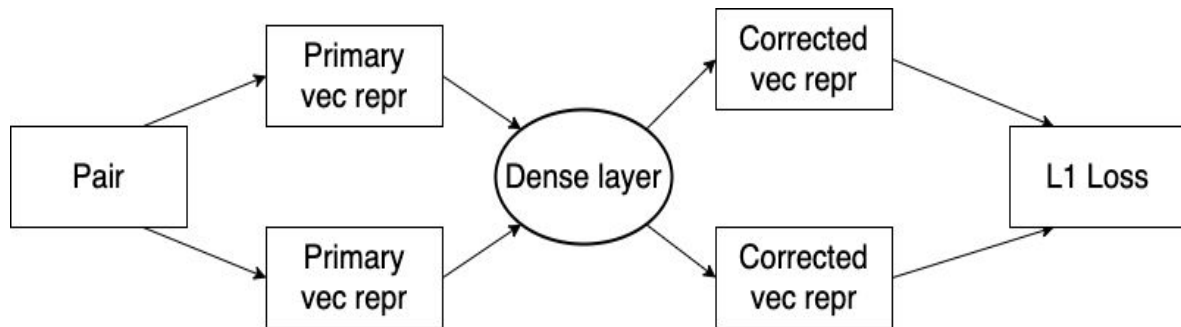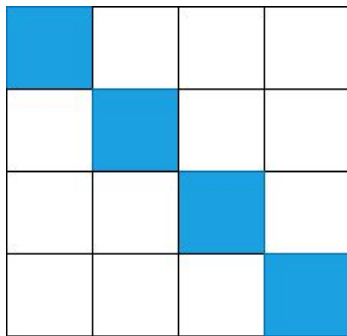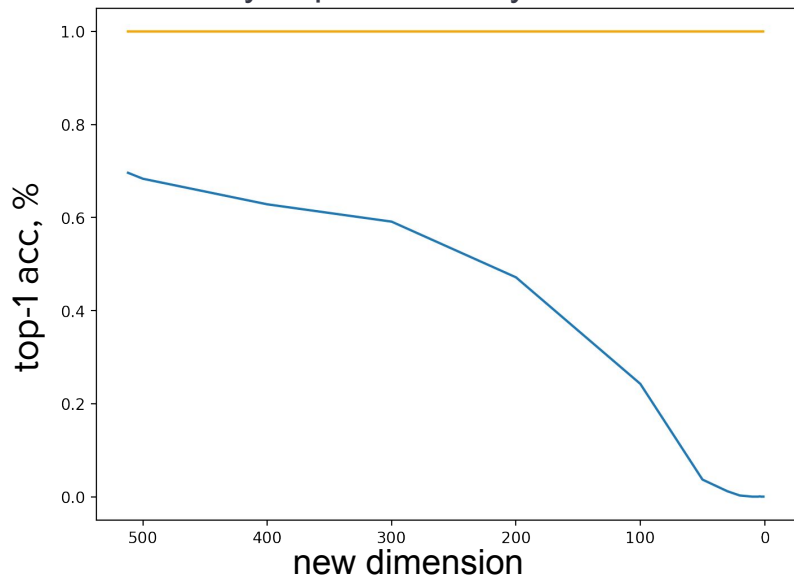
# CLIP basics

# Dimensionality Reduction

Alexey Kolosov, Ekaterina Orlova

**Task**: Investigate dimensionality reduction methods and show that for the presented data there exist such embedding dimensionality D' < D which doesn't decrease embeddings correspondence quality (top-1 accuracy).

Quality dependence by new dimension

# **Neural MDS**

Alexey Kolosov

**Problems statements**

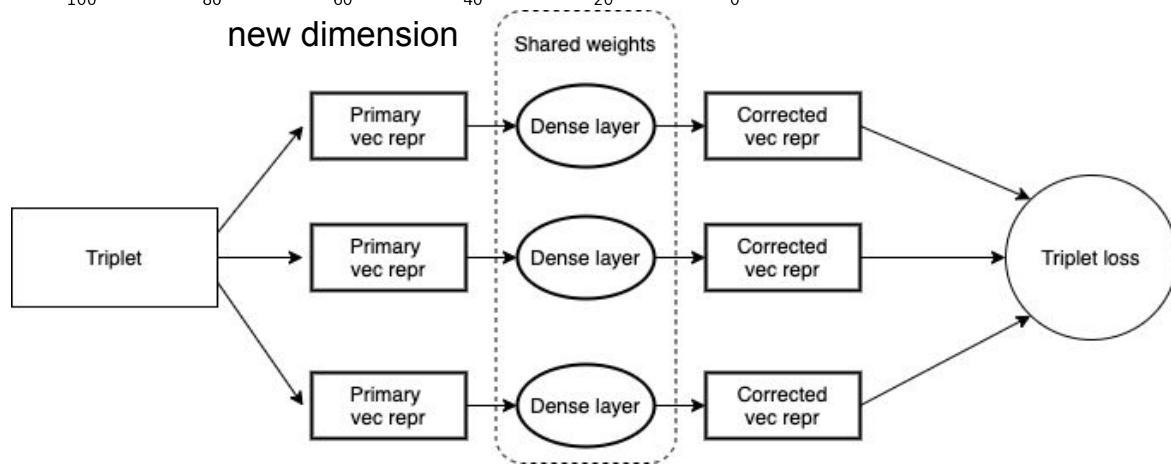1. **COCO, isometric, val**
2. **COCO, isotonic, val**

$$d(i, j) = e(g(i), g(j))$$
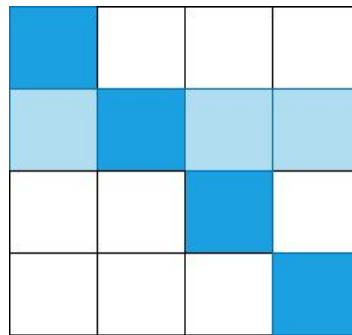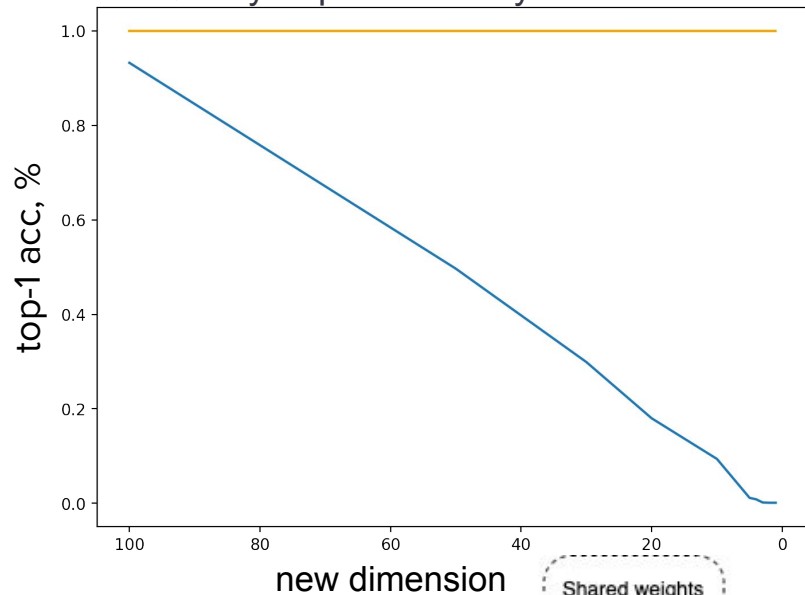
Top-1 accuracy

COCO, isometric, val

5000 pairs, 200 epochs

for 512 dim - 1639 pairs

6

Quality dependence by new dimension

top-1 acc, % (y-axis)
new dimension (x-axis)

# Neural MDS

Alexey Kolosov

**Isotonic problem result**

$$d(i,j) < d(k,l) \Rightarrow$$
$$e(g(i), g(j)) < e(g(k), g(l))$$

Top-1 accuracy

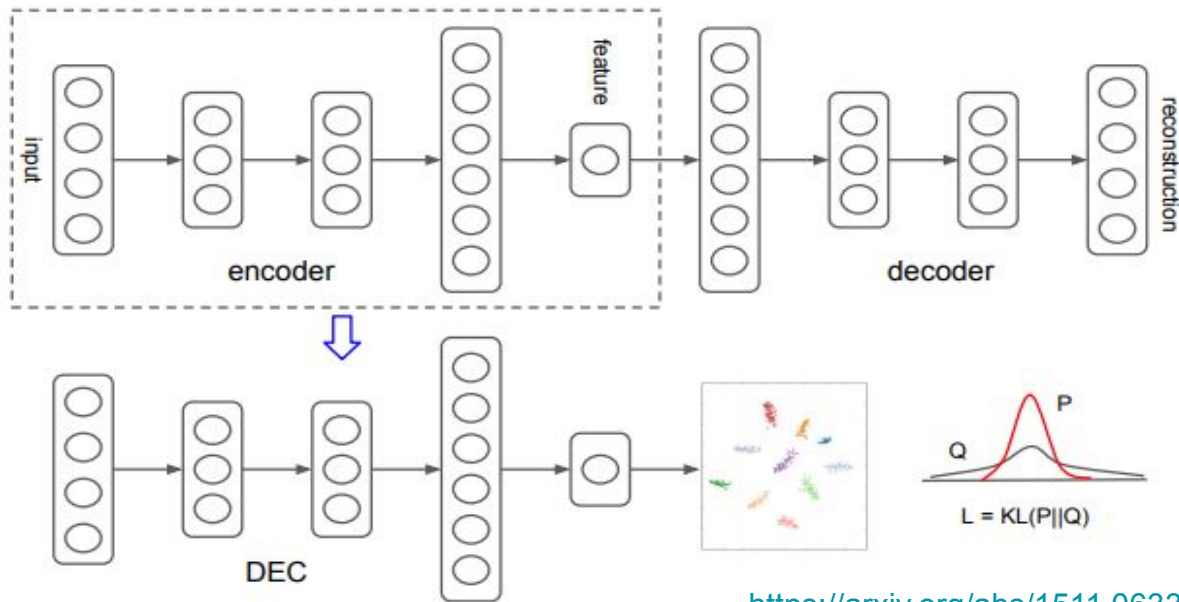COCO, isotonic, val

8000000 quads, 200 epochs

for 512 dim - 1639 pairs

7

# Embeddings clusterization

Abdullaeva Uma, Anna Dmitrienko, Sergey Skorik, Anna Rudenko

**Task**: Investigate the clusterization methods and show that there exist clusters in the embeddings data. Perform the visualization of these data clusters.

**Deep embedded clustering (DEC) model**



https://arxiv.org/abs/1511.06335

8

Let "S"- set of n element
$X = \{X1, X2, \ldots, Xn\}$ — the division into classes
$Y = \{Y1, Y2, \ldots, Yn\}$ - the resulting division into clusters

| $X \backslash Y$ | $Y_1$ | $Y_2$ | $\ldots$ | $Y_s$ | Sums |
|---|---|---|---|---|---|
| $X_1$ | $n_{11}$ | $n_{12}$ | $\ldots$ | $n_{1s}$ | $a_1$ |
| $X_2$ | $n_{21}$ | $n_{22}$ | $\ldots$ | $n_{2s}$ | $a_2$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ | $\vdots$ |
| $X_r$ | $n_{r1}$ | $n_{r2}$ | $\ldots$ | $n_{rs}$ | $a_r$ |
| Sums | $b_1$ | $b_2$ | $\ldots$ | $b_s$ | $n$ |

$$p_{ij} = \frac{n_{ij}}{n}, p_i = \frac{a_i}{n}, p_j = \frac{b_j}{n}$$

# Metrics

$$\underbrace{\underbrace{ARI}_{\text{Adjusted Index}}} = \frac{\overbrace{\sum_{ij}\binom{n_{ij}}{2}}^{\text{Index}} - \overbrace{[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}]/\binom{n}{2}}^{\text{Expected Index}}}{\underbrace{\frac{1}{2}[\sum_i \binom{a_i}{2}+\sum_j \binom{b_j}{2}]}_{\text{Max Index}} - \underbrace{[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}]/\binom{n}{2}}_{\text{Expected Index}}}$$

$$\mathrm{ARI} = \frac{\mathrm{RI} - E[\mathrm{RI}]}{\max(\mathrm{RI}) - E[\mathrm{RI}]} \qquad \mathrm{RI} = \frac{a+b}{C_2^{n_{samples}}}$$

$$\mathrm{NMI}(U,V) = \frac{\mathrm{MI}(U,V)}{\mathrm{mean}(H(U), H(V))}$$

# **Results**

|  | AMI | ARI | FMI | NMI |
|---|---|---|---|---|
| K-Means | 0.72 | 0.63 | 0.67 | 0.72 |
| Auto-encoder | 0.74 | 0.67 | 0.70 | 0.76 |
| **DEC** | **0.82** | **0.77** | **0.79** | **0.82** |

Visualizing clusters using T-SNE

# Benchmark
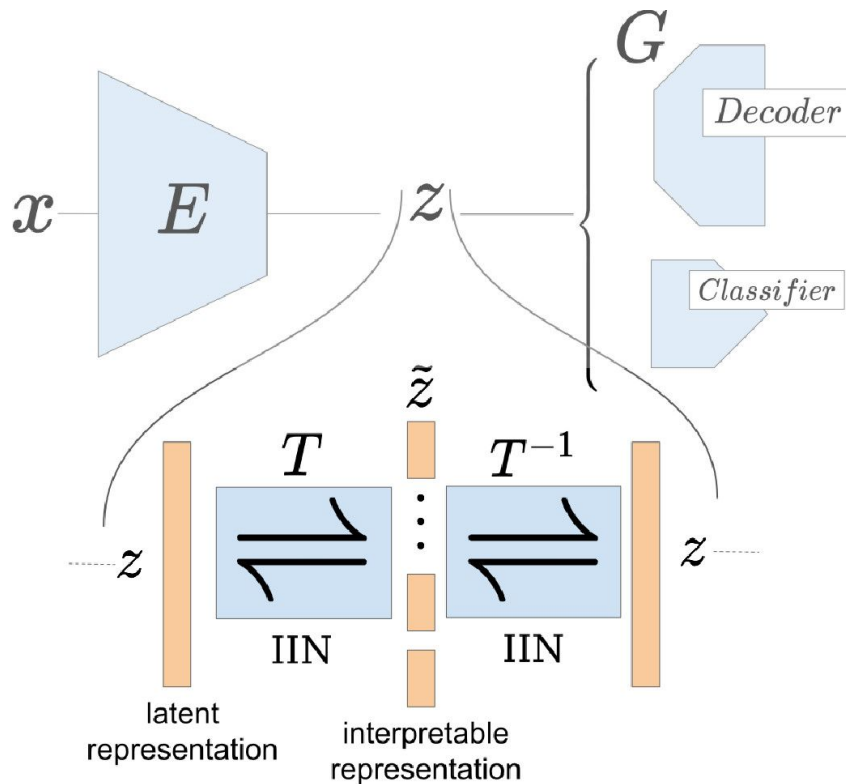


Image Clustering on CIFAR-10

# Advantages of DEC model

- DEC method is linear in the number of data points and scales gracefully to large datasets

- DEC employs deep neural networks to perform non-linear embedding that is necessary for more complex data

- CLIP + DEC show SOTA results in clustering

# Disentanglement

Ekaterina Orlova, Anna Dmitrienko, Sergey Skorik, Anna Rudenko, Abdullaeva Uma



Embedding in latent space:
$$z = E(x) \in \mathbb{R}^{H \times W \times C}$$

Invertible Interpretation Network:
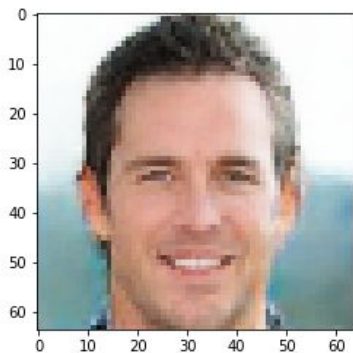$$T(z) = \bar{z}$$

Modified latent vector z:
$$z \to z^\star := T^{-1}(T(z)^\star)$$

Loss function:
$$\mathcal{L} = \sum_{F=1}^{K} \mathbb{E}_{(x^a, x^b) \sim p(x^a, x^b | F)} \ell\left(E(x^a), E(x^b) \mid F\right),$$
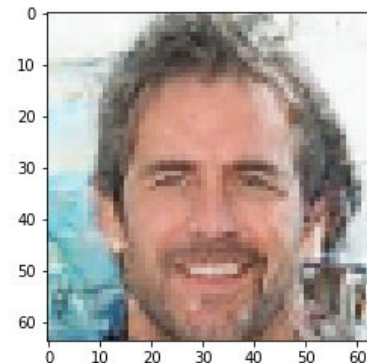
where $l$ — per-example loss

13

# Process



| | |
|---|---|
| -1,99 | -1,99 |
| -1,65 | -1,65 |
| 1,94 | 1,94 |
| -0,82 | -0,82 |
| -0,82 | -0,82 |
| -1,06 | -1,06 |
| -1,23 | -1,23 |
| 1,06 | 1,06 |
| 0,32 | 0,32 |
| 2,52 | 2,52 |
| -0,28 | -0,28 |
| -0,40 | -3,00 |
| 0,18 | 0,18 |

$E$
$T$

$D$
$T-1$

# SelebA - Glasses



original

2 component = 1.2

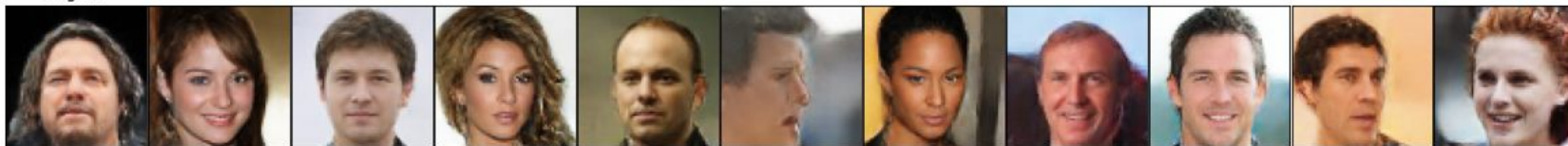2 component = 2

# SelebA - Race



original
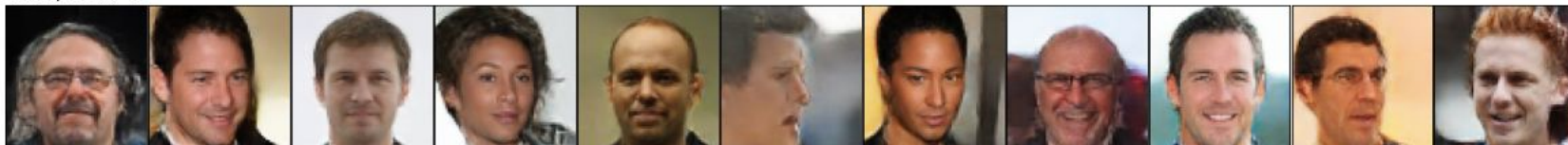
5 component = -1

5 component = 3

# SelebA - Sex



original

1 component = -3

1 component = 2

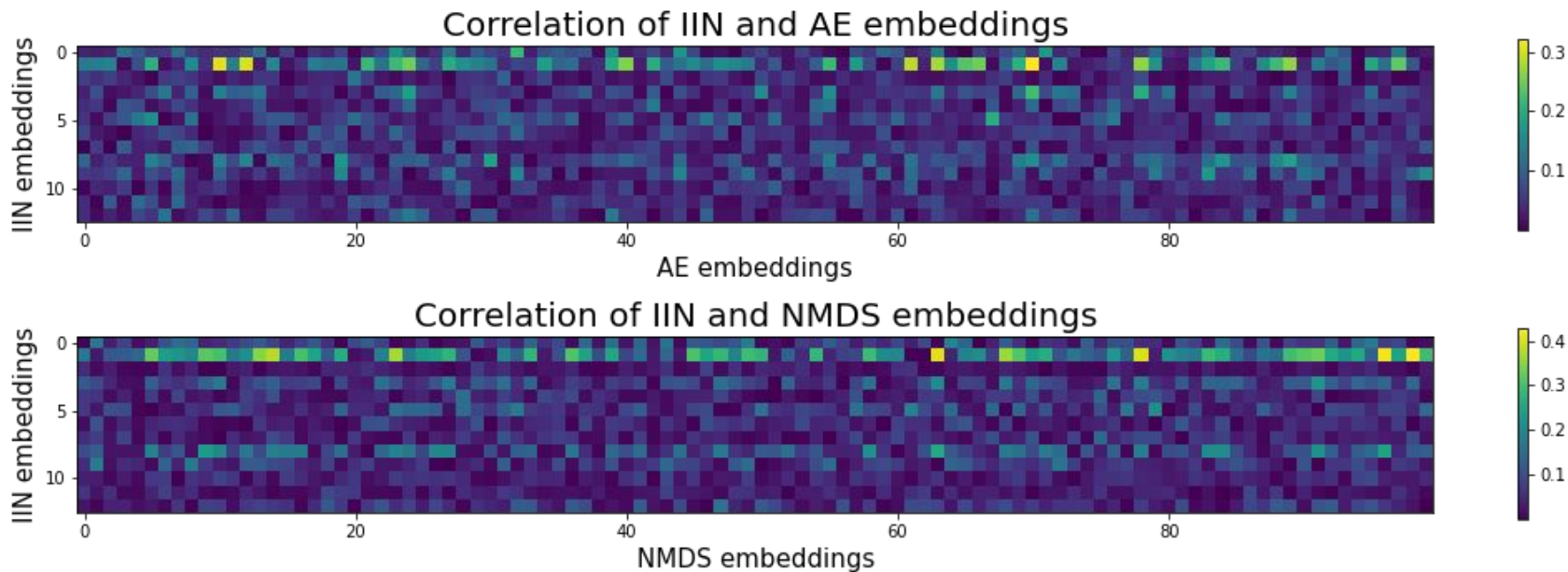# Semantic analysis by CLIP

## Cosine similarity between text and image features



|  | | | | | |
|---|---|---|---|---|---|
| **woman with red hair** | 0.2810 | 0.2089 | 0.2239 | 0.2721 | 0.2367 | 0.1793 | **man with hat and sunglasses** |

(table continues below with row labels)

| Row label (left) | img1 | img2 | img3 | img4 | img5 | img6 | Row label (right) |
|---|---|---|---|---|---|---|---|
| woman with red hair | 0.2810 | 0.2089 | 0.2239 | 0.2721 | 0.2367 | 0.1793 | man with hat and sunglasses |
| blonde woman with red lipstick | 0.2429 | 0.2723 | 0.2194 | 0.2269 | 0.2520 | 0.1774 | man with glasses turned his face |
| woman with sunglasses | 0.2659 | 0.2224 | 0.2888 | 0.2005 | 0.2182 | 0.2450 | woman with long hair |

# Correlation of semantic features and embeddings



Correlation of IIN and AE embeddings

Correlation of IIN and NMDS embeddings

# **Conclusions**

- Proposed new method for **dimensionality reduction** of CLIP embeddings
- Proposed and evaluated new method for **clusterization**, close to SOTA
- Interpreted **latent representations** of various VAE

# What's next?

- Automatic data augmentation with text descriptions
- Improving managed image synthesis