

# Galaxy Zoo 2 - galaxy clusterization

Autorzy:

**Paweł Pozorski**

**Michał Pytel**

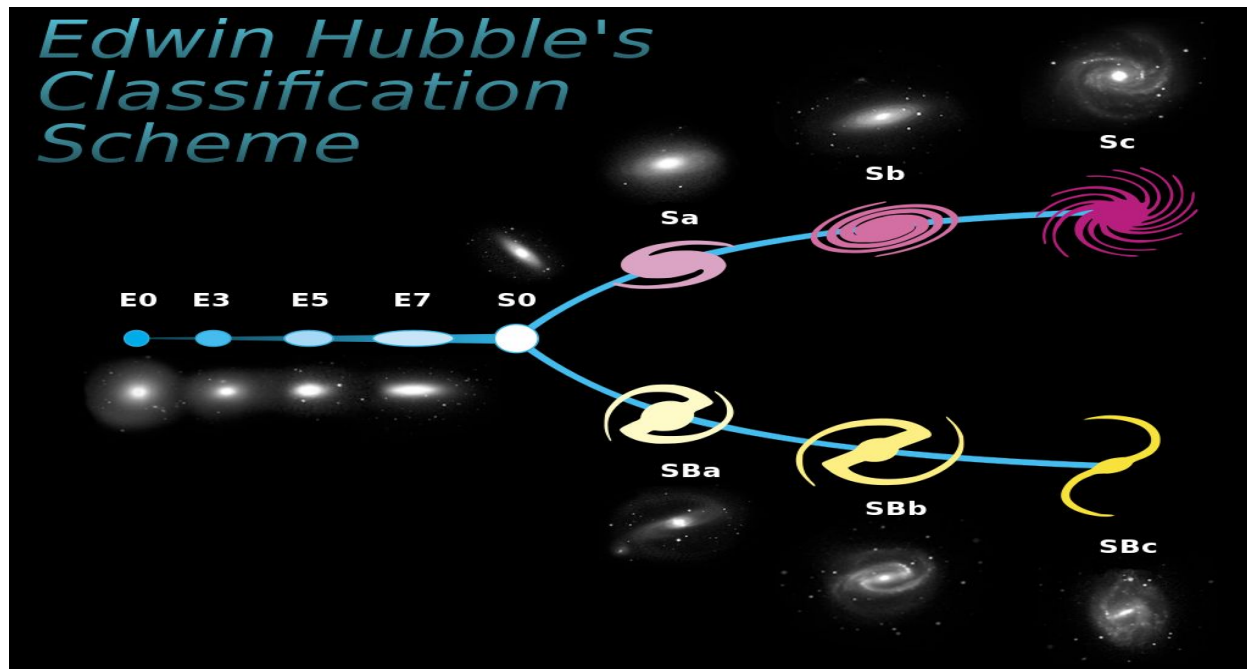
**Michał Matuszyk**

# Troszkę o danych

- Źródło:  
[www.kaggle.com/datasets/jaimetrickz/galaxy-zoo-2-images](https://www.kaggle.com/datasets/jaimetrickz/galaxy-zoo-2-images)
- Ponad 300k kolorowych zdjęć (355990)
- (424x424x3)

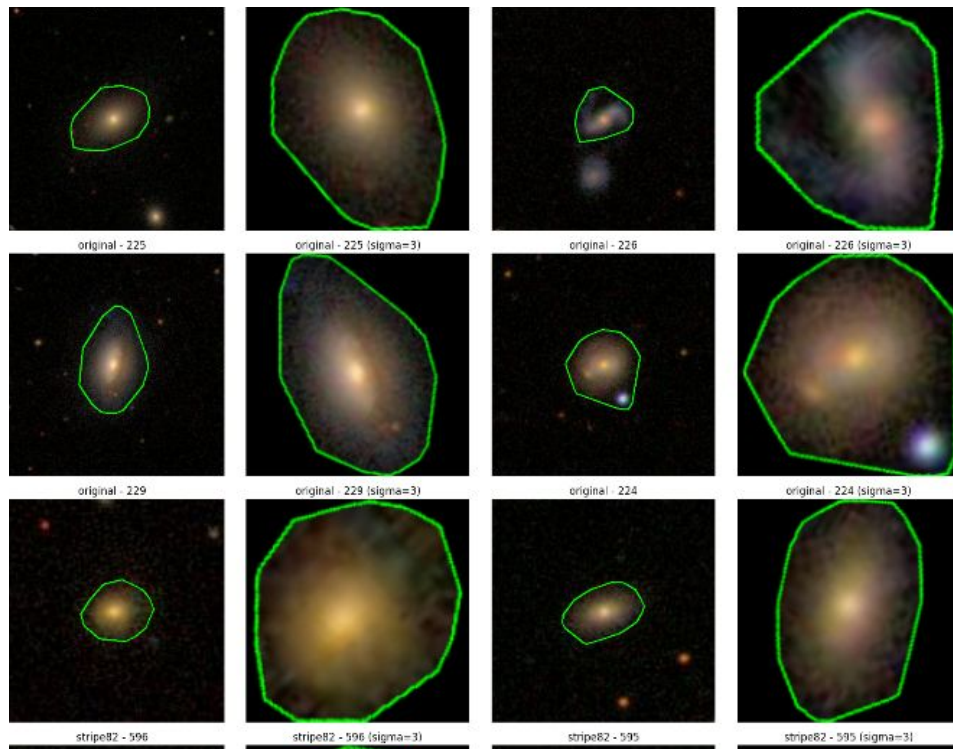


# Cel biznesowy

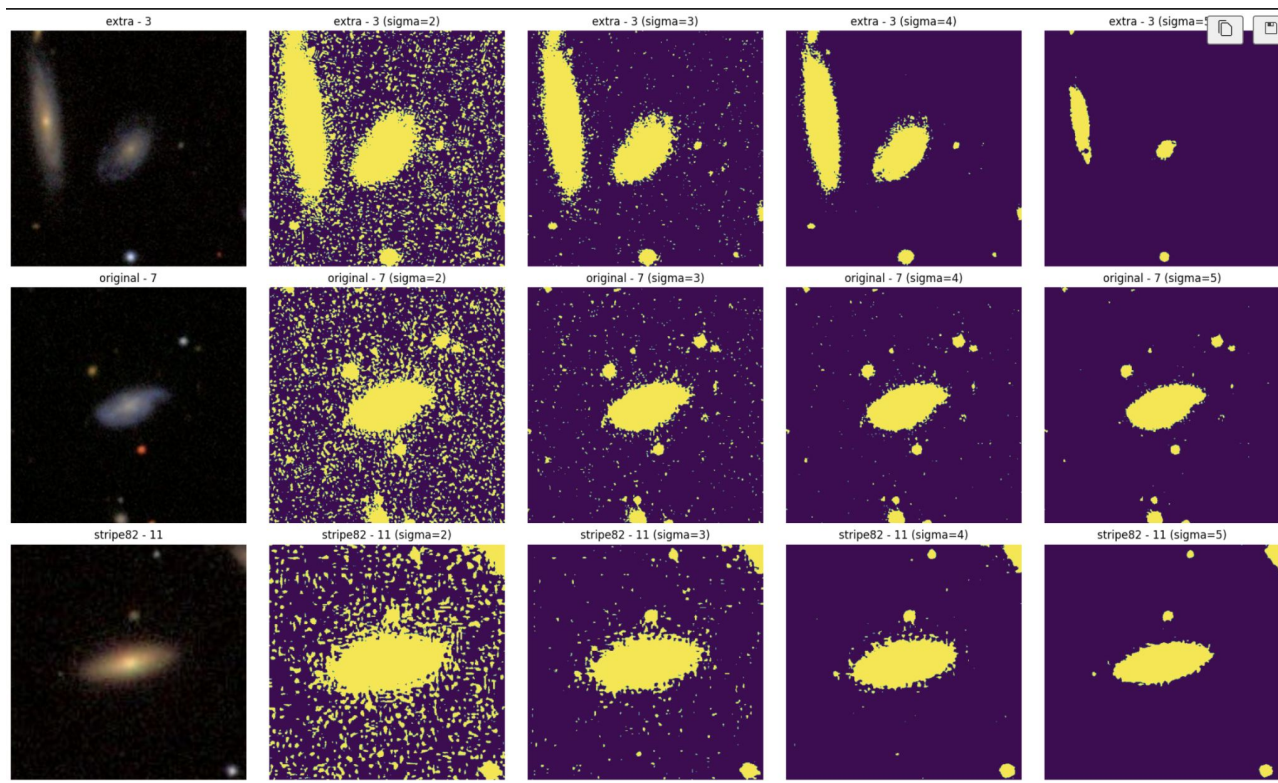


Sprawdzenie czy da się wykazać zachodzenia sekwencji Hubble'a

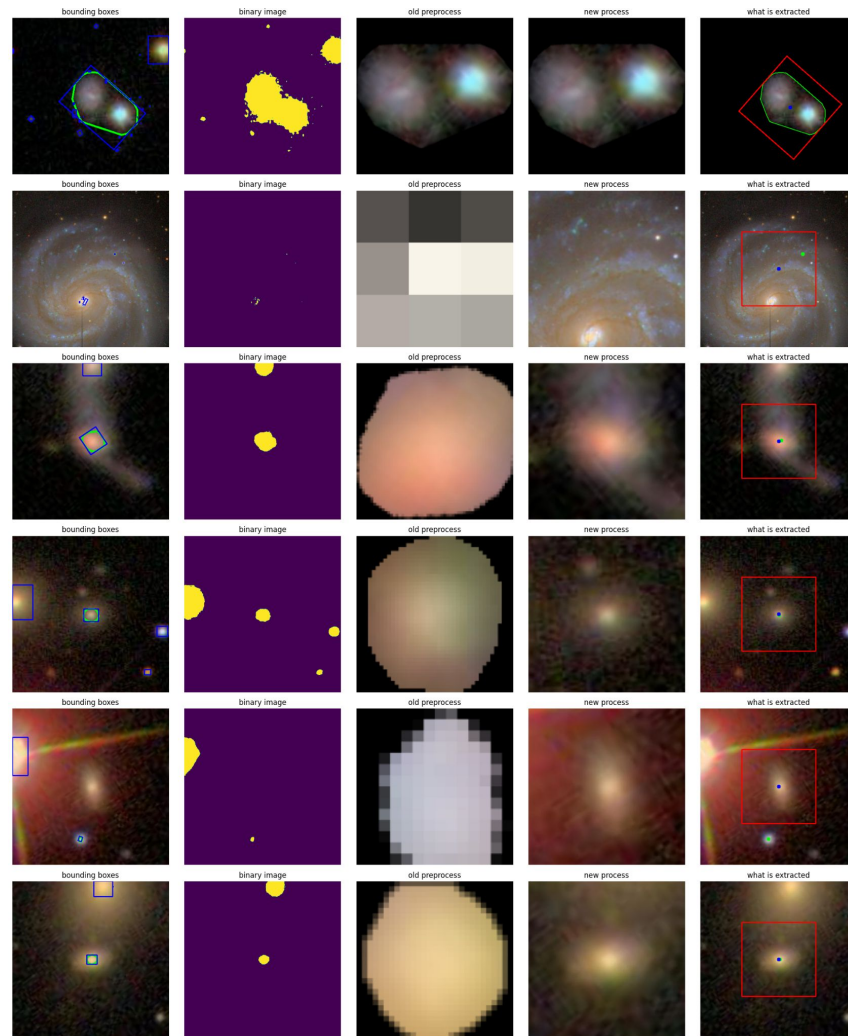
# Proces oczyszczania danych



# Proces oczyszczania danych



# Proces oczyszczania danych



# Proces oczyszczania danych

Empirycznie - proces czyszczenia zdjęcia Z0 składa się z kolejnych kroków:

1. Zdjęcie jest rzucane do GrayScale Z1
2. Na podstawie Z1 obliczane jest mean, median i std
3. Z1 jest rzucane do binarnego Z2 po przeczyszczeniu -  $\leq \text{median} + (\text{sigma} * \text{std})$  jest ustawiana jako 0. Przyjęta finalnie sigma to 4
4. Na podstawie Z2 szukamy konturów obiektów, z wszystkich znalezionych konturów wybieramy K maksymalizujący  
`cv2.contourArea(c)-contour_distance_from_center(c) ** 2`
5. K1 służy jako maska dla Z0 - wszystko poza jego wnętrzem jest ustawiane na 0 - dostajemy Z3
6. Z Z3 wycinany jest minimalny prostokąt zawierający non-0 pixels a następnie rozszerzany i skalowany do kwadratu 150\*150 z zachowaniem aspektu

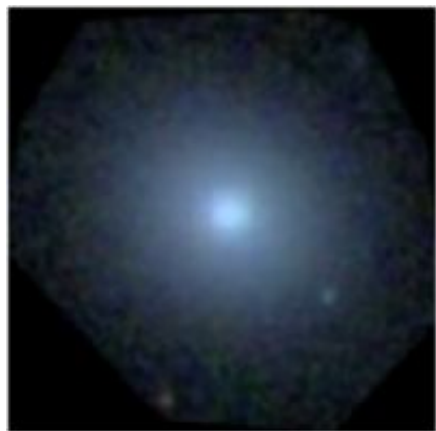
! UWAGA - przeprocesowane zdjęcia sprawiają wrażenia niebieskich ponieważ cv2 zmienia warstwę R z B.

# Podejście 1

Sam autoencoder jako compressor i feature extractor -> KMeans

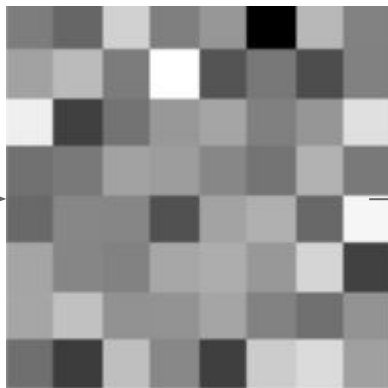


# Model



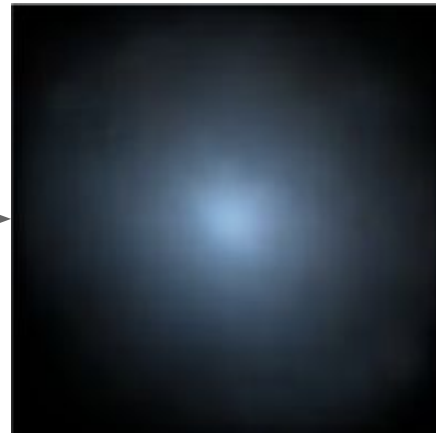
Przeczyszczone  
zdjęcie 150x150x3

Encoding



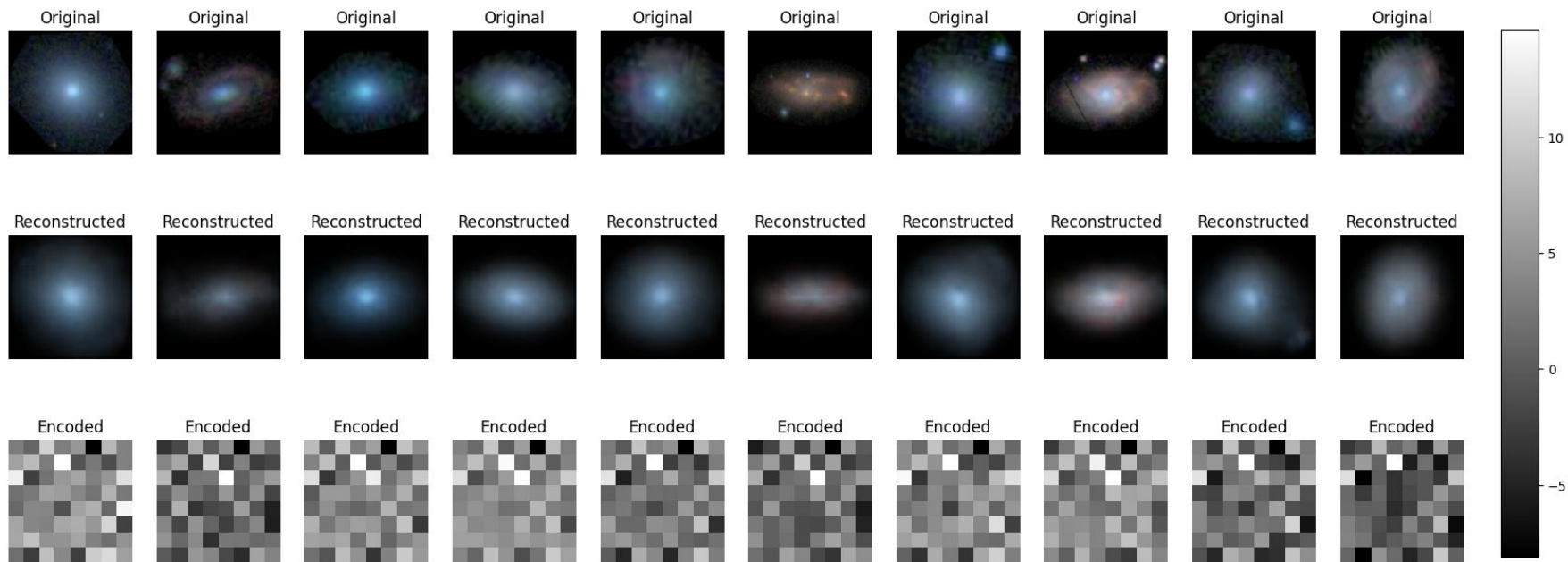
Cechy  
Wymiar:

Decoding



Poglądowo -  
czego nasz model  
się nauczył?

# Efekty



# Podejście 2

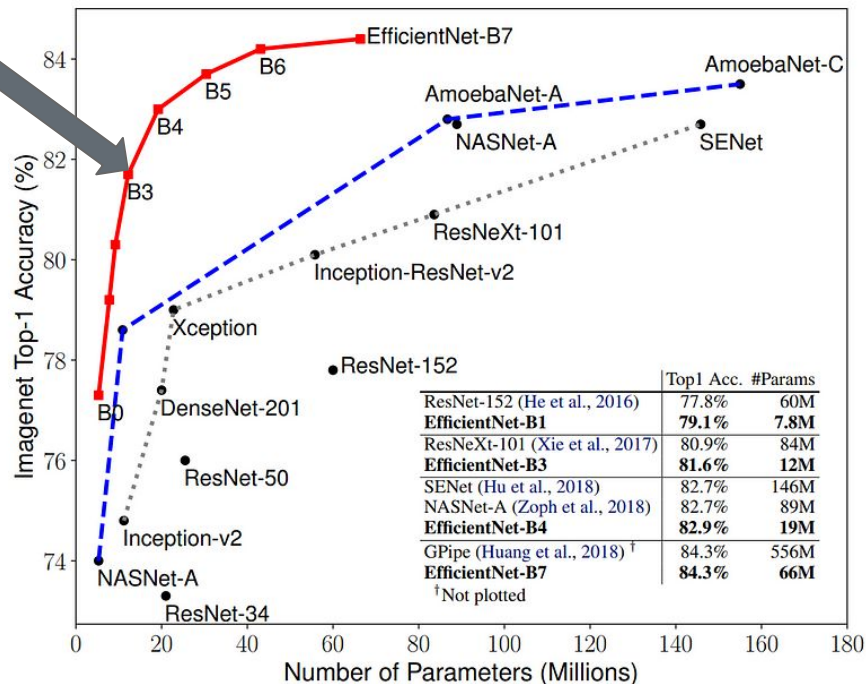
Warstwy konwolucyjne z EfficientNet B3 z ImageNetu -> autoencoder  
jako compressor i feature extractor -> PCA -> clustering

EfficientNet zwraca wektor  $4 \times 4 \times 1536$

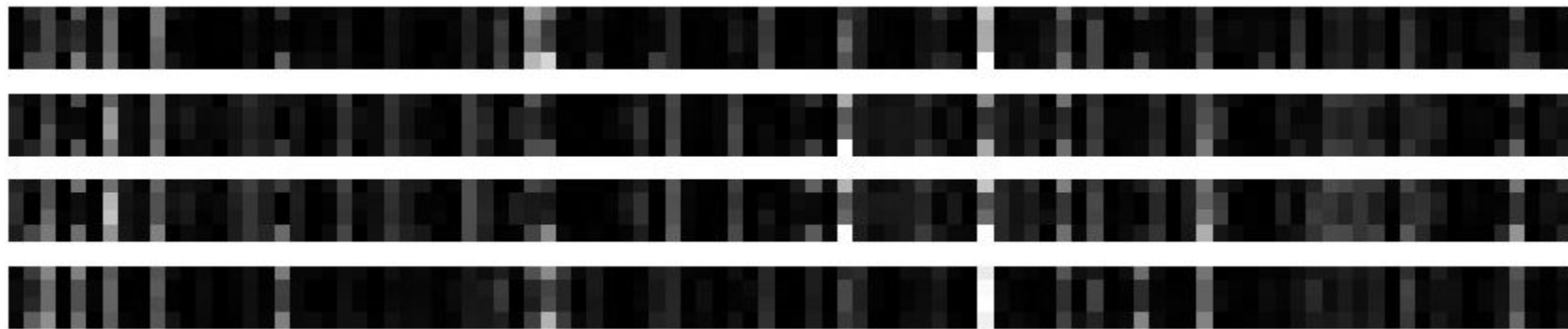
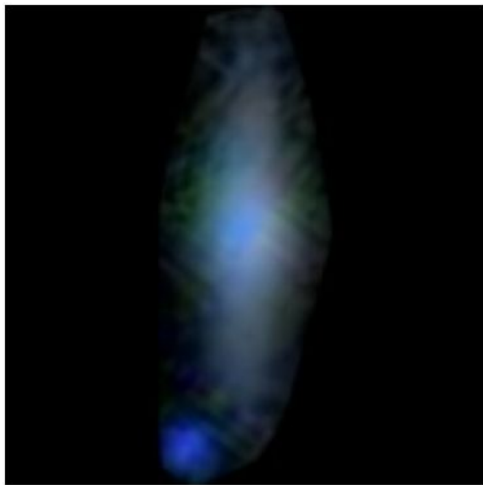
Autoencoder wektor  $1 \times 256$

PCA wektor  $1 \times 8$  (zachowanie wariancji: 0.99)

# Wybór Feature Extactora

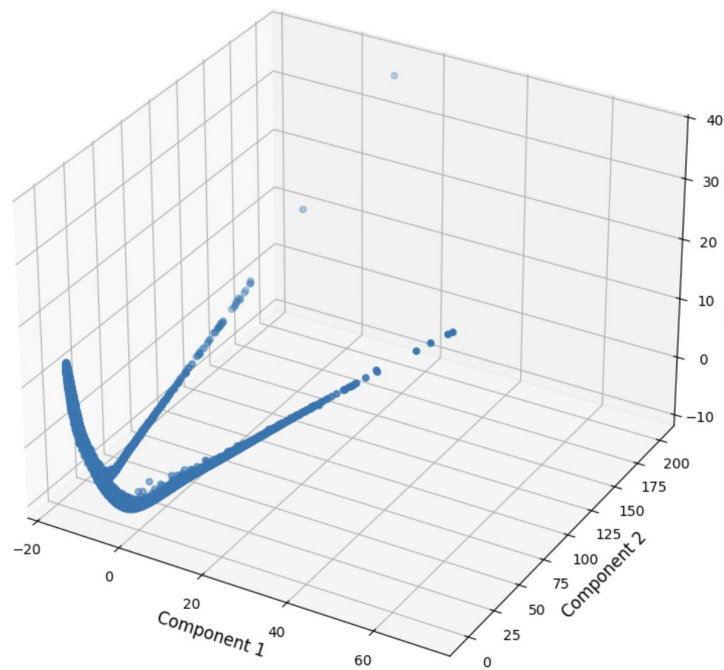


# EfficientNet Features

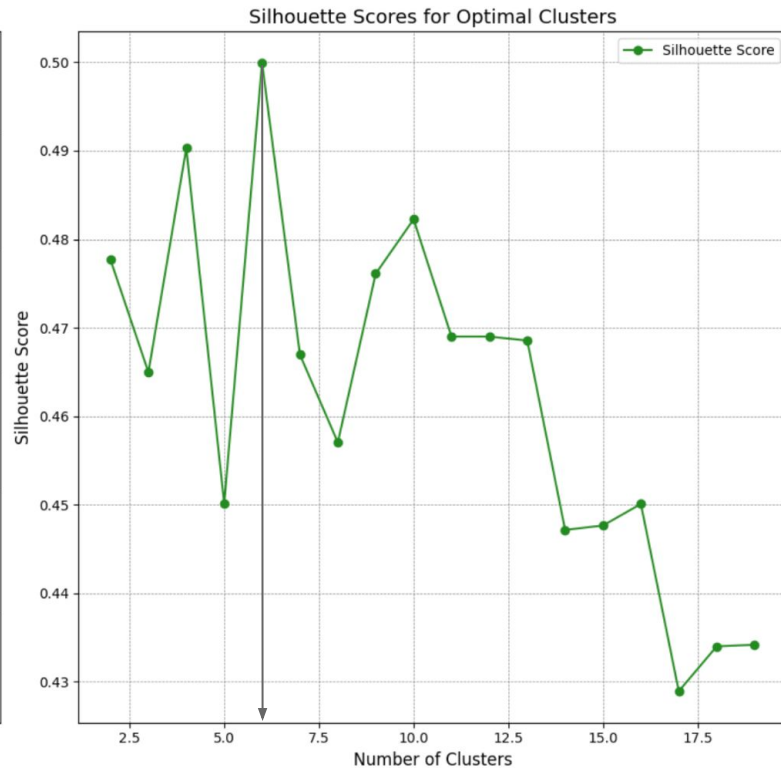
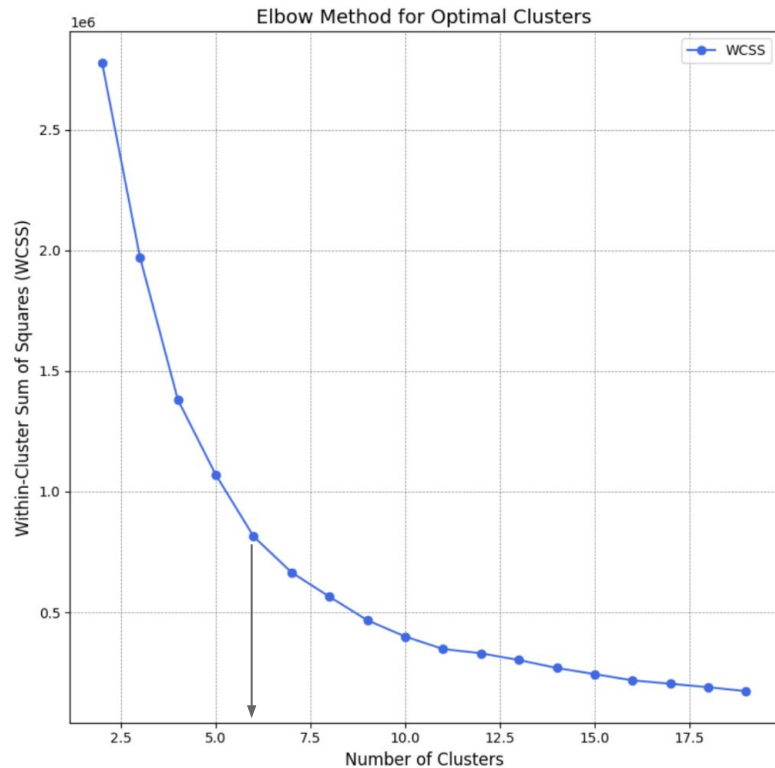


# PCA

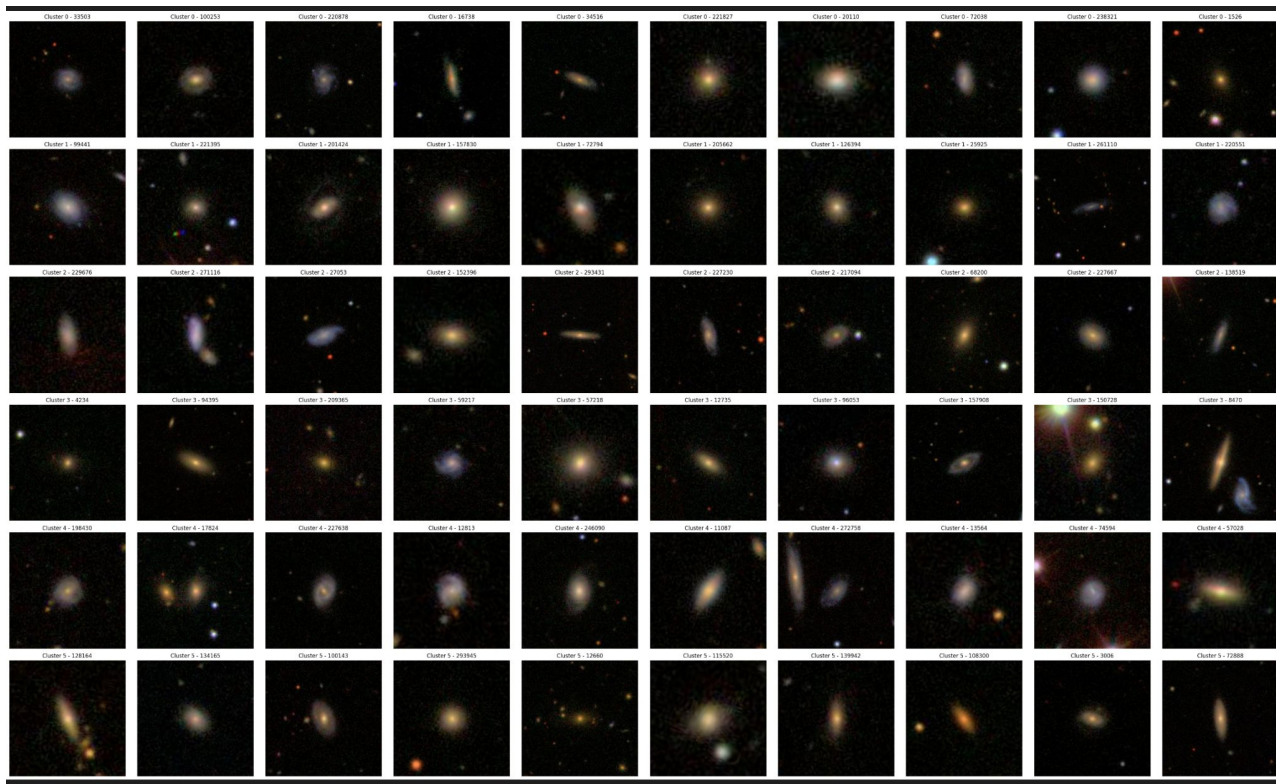
PCA - First 3 Components



# KMeans – wybór liczby klastrów

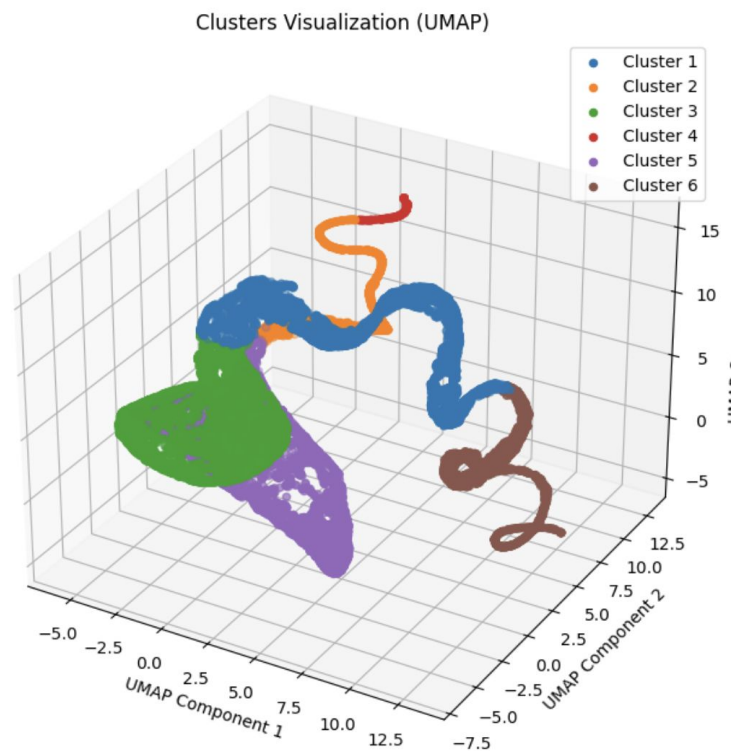
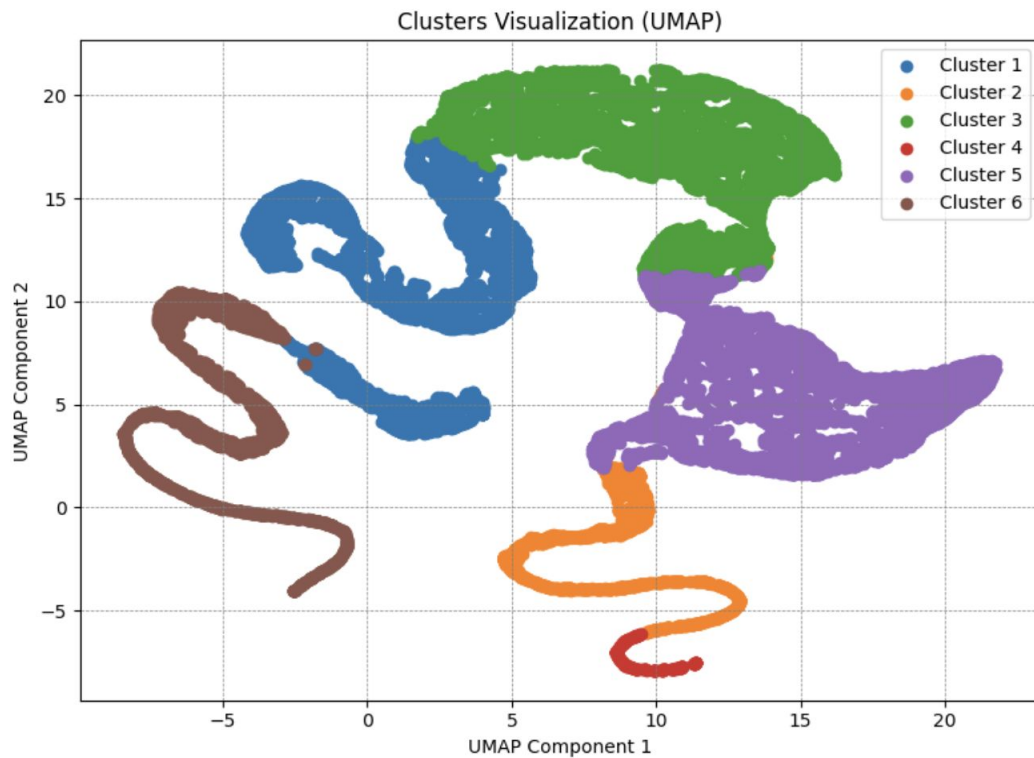


# KMeans – Efekty

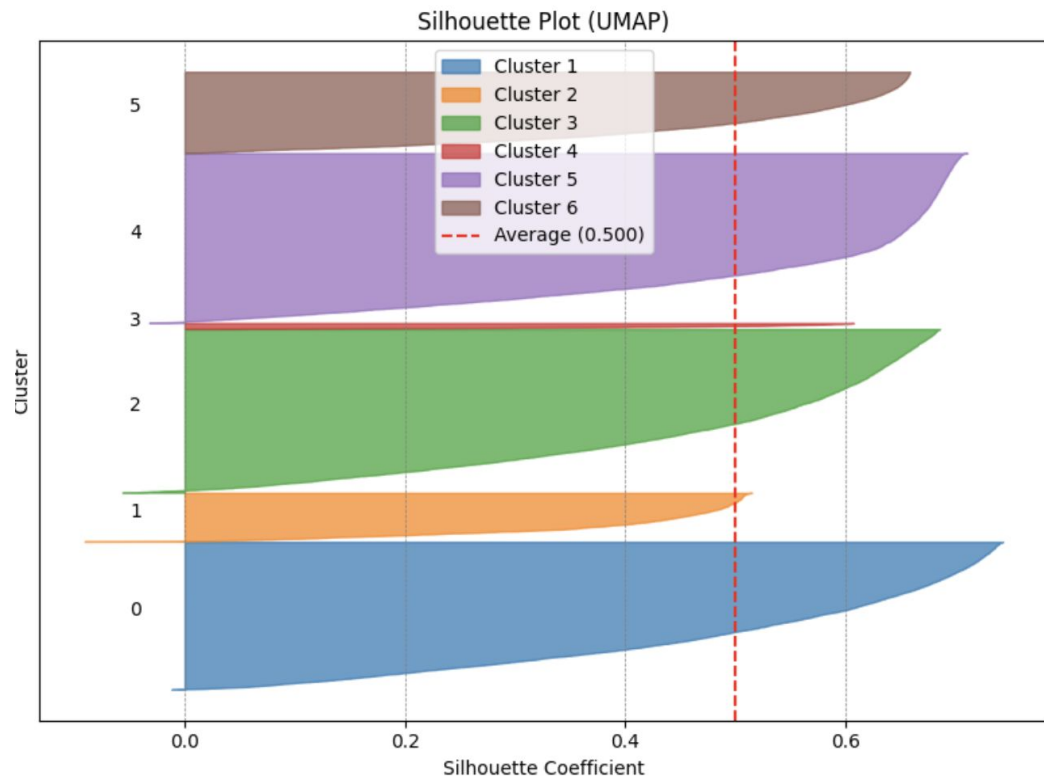




# KMeans – Efekty – UMAP 2D, 3D



# KMeans – Efekty



# Podejście 3

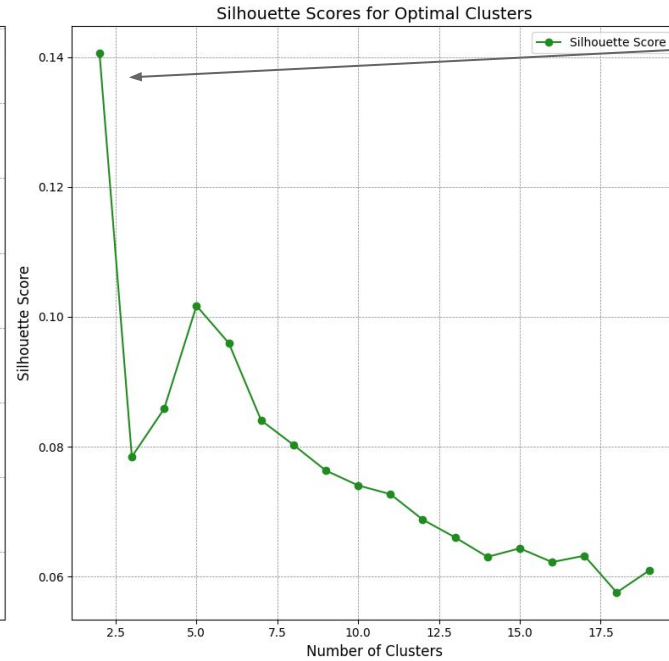
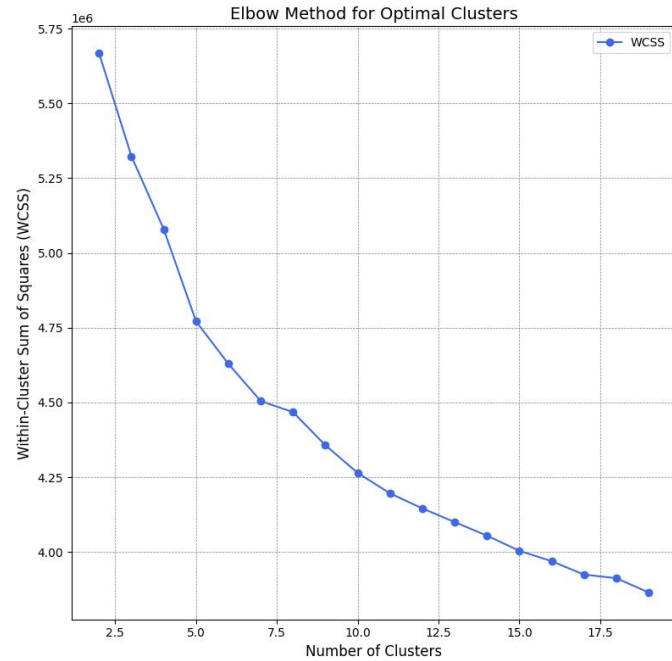
preprocesowanie danych -> feature extractor - WGG16 -> PCA -> clustering

WGG16: zwraca wektor 1x4096

PCA: zwraca 1x1387 (zachowuje 99% wariancji)

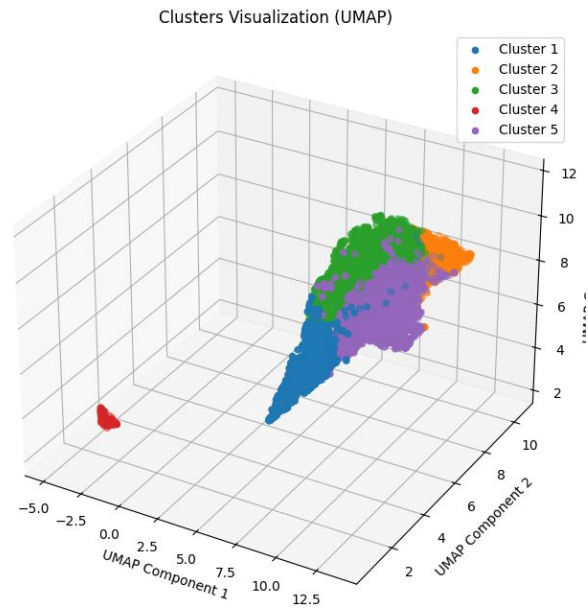
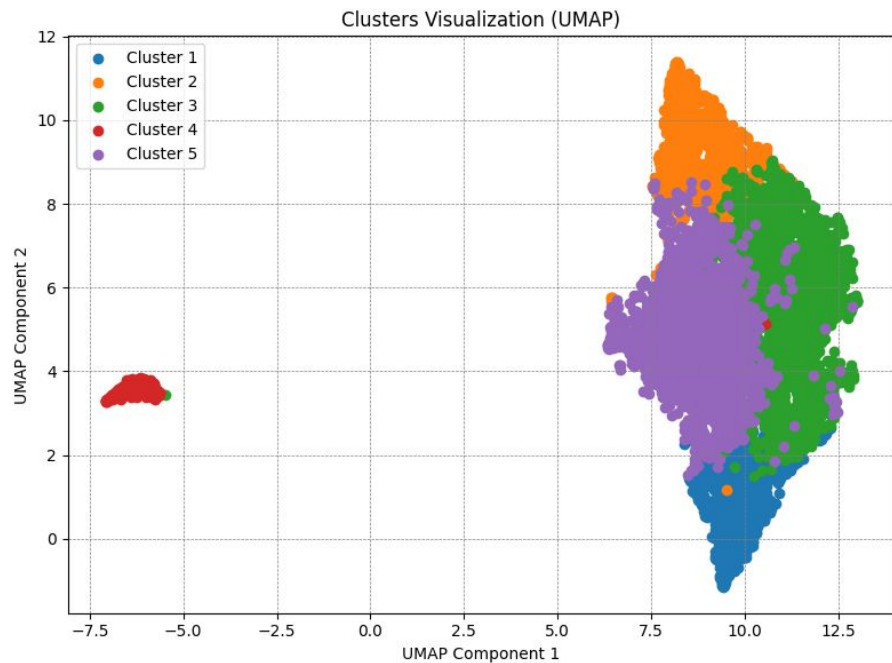
Clustering: 5 Klastrów

# KMeans i silhouette



Low score

# Wyjaśnialność



# Końcowy wybór

- Wybieramy podejście 2
- Najlepszy Silhouette score: 0,52
- Najlepsza interpretowalność
- Indeks Daviesa-Bouldina: 0.7

Dziękujemy za  
uwagę

