

# CNNによるスタイル変換と Web画像を用いた 画像の任意質感生成

松尾 真, 柳井 啓司

電気通信大学 大学院情報理工学研究科  
総合情報学専攻

はじめに

## Neural Style Transfer

➤ Deep Neural Networkを用いた画像変換

➤ 写真と絵画を入力すると、写真が絵画調に

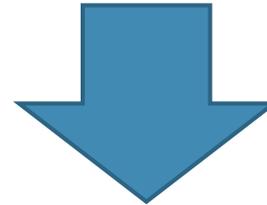
➤ DNNによる画像合成分野が発足



content image



style image



generated image

# 目的

## 単語概念からの自動スタイル生成

画像内物体の質感の任意変換への応用

- ▶ 質感別画像データの増量
- ▶ デザイン, エンターテインメント  
様々な分野への応用が期待できる



content  
image

「金属製」



「木製」



request

generated image

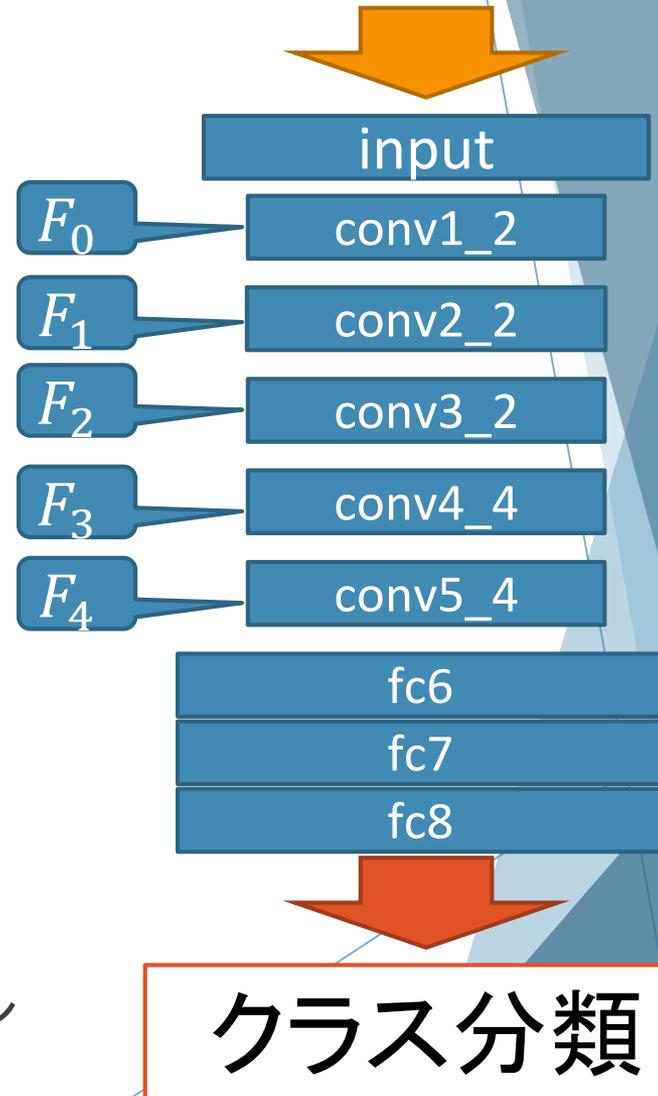
# VGG19 Network (Max pool層等省略)

## 関連研究

### Neural Style Transfer

L. A. Gatys et al, “A neural algorithm of artistic style,” in arXiv:1508.06576, 2015.  
[Gatys et al.]

- ImageNet1000クラスの分類を学習済みのDNN
- 画像入力時に各Layer内でフィルタの活性値行列 $F_i$ が発生する

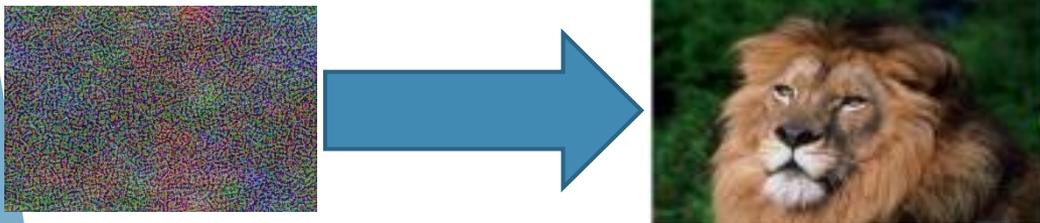


# VGG19 Network (Max pool層等省略)

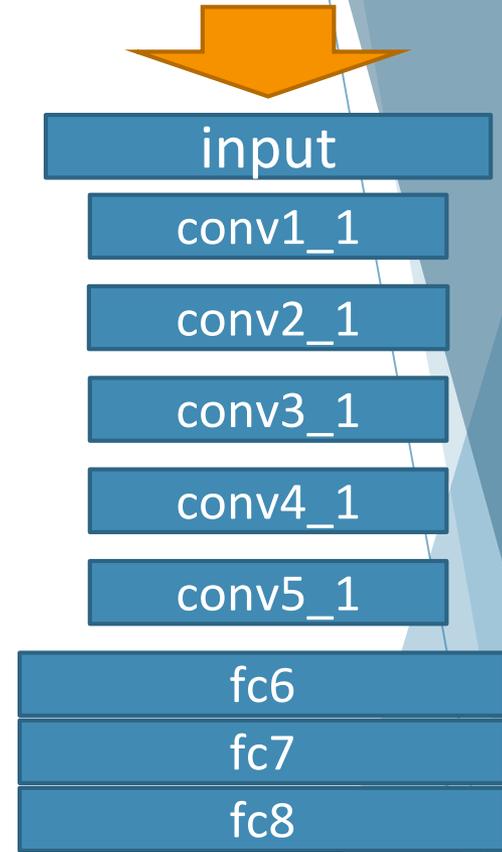
## 関連研究

Neural Style Transfer  
[Gatys et al.]

生成画像と元画像の  $F_l$  のずれを最小化するように生成画像を最適化



元画像を復元！

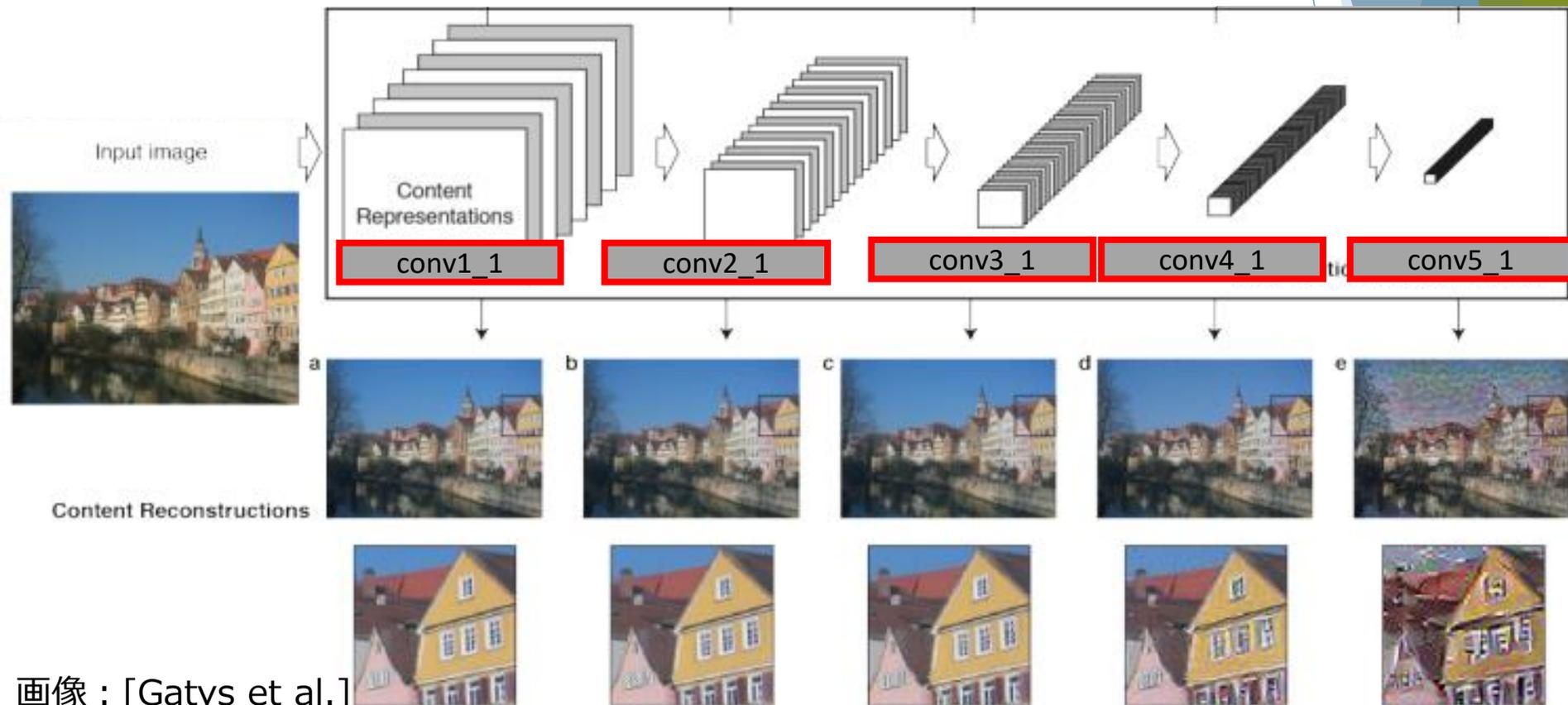


Classification

# Neural Style Transfer

[Gatys et al.]

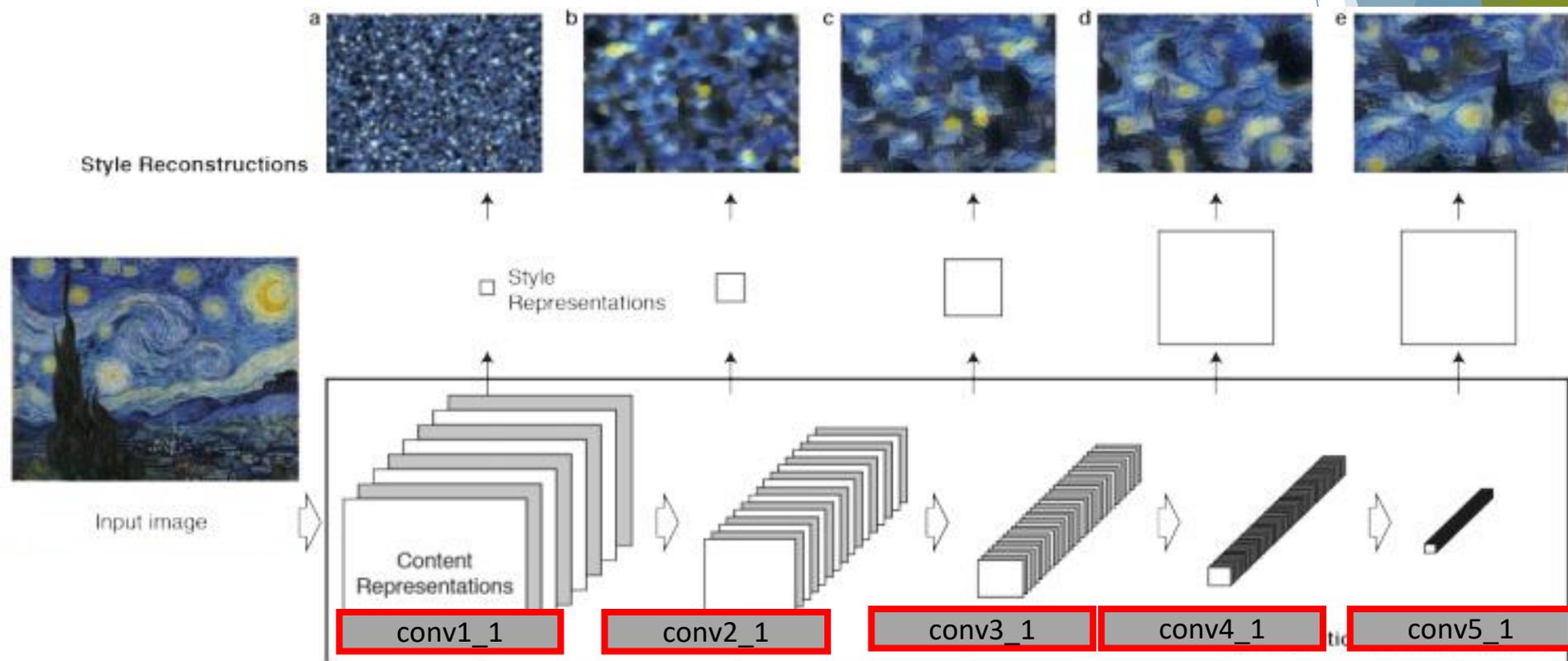
コンテンツ表現： $F_l$ による画像の復元



# Neural Style Transfer

[Gatys et al.]

スタイル表現 :  $G_l = F_l F_l^T$  による画像の復元



# Neural Style Transfer

## [Gatys et al.]

二種の画像表現から損失関数を形成

- ▶ **コンテンツ表現**：深い層における  $F_l$ 
  - ▶ 形状のみを保持
- ▶ **スタイル表現**：様々な層における  $G_l = F_l F_l^T$ 
  - ▶ スタイルのみを保持
- ▶ 2枚の異なる画像の形状, スタイルを持った画像を復元

入力



コンテンツ画像  $x_c$



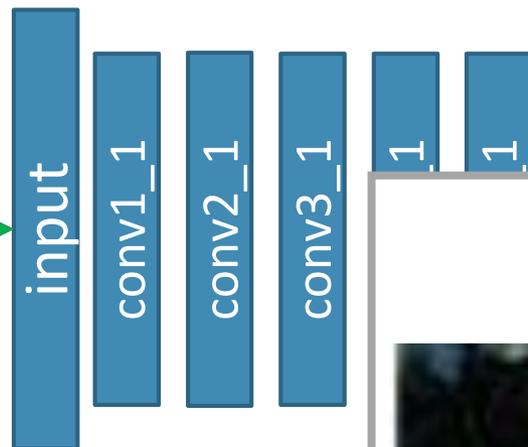
スタイル画像  $x_s$

出力 (初期状態)



生成画像  $x_g$

Pre-trained  
DNN(VGG19)



$F(x_c)$

$F(x_g)$

$Loss_{content}(x_c, x_g)$

出力 (最終状態)



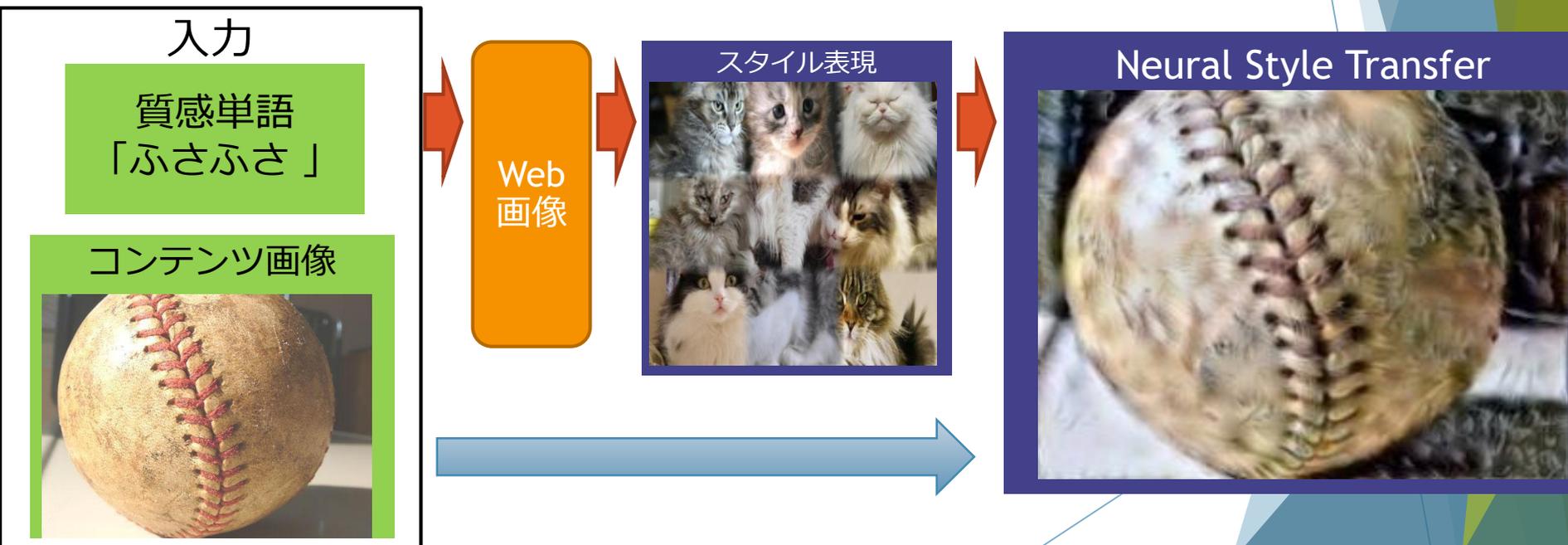
生成画像  $x_g$

画像最適  
(L-BFGS)

生成画像のコンテンツをコンテンツ画像に，スタイルをスタイル画像に近づける

# 提案手法

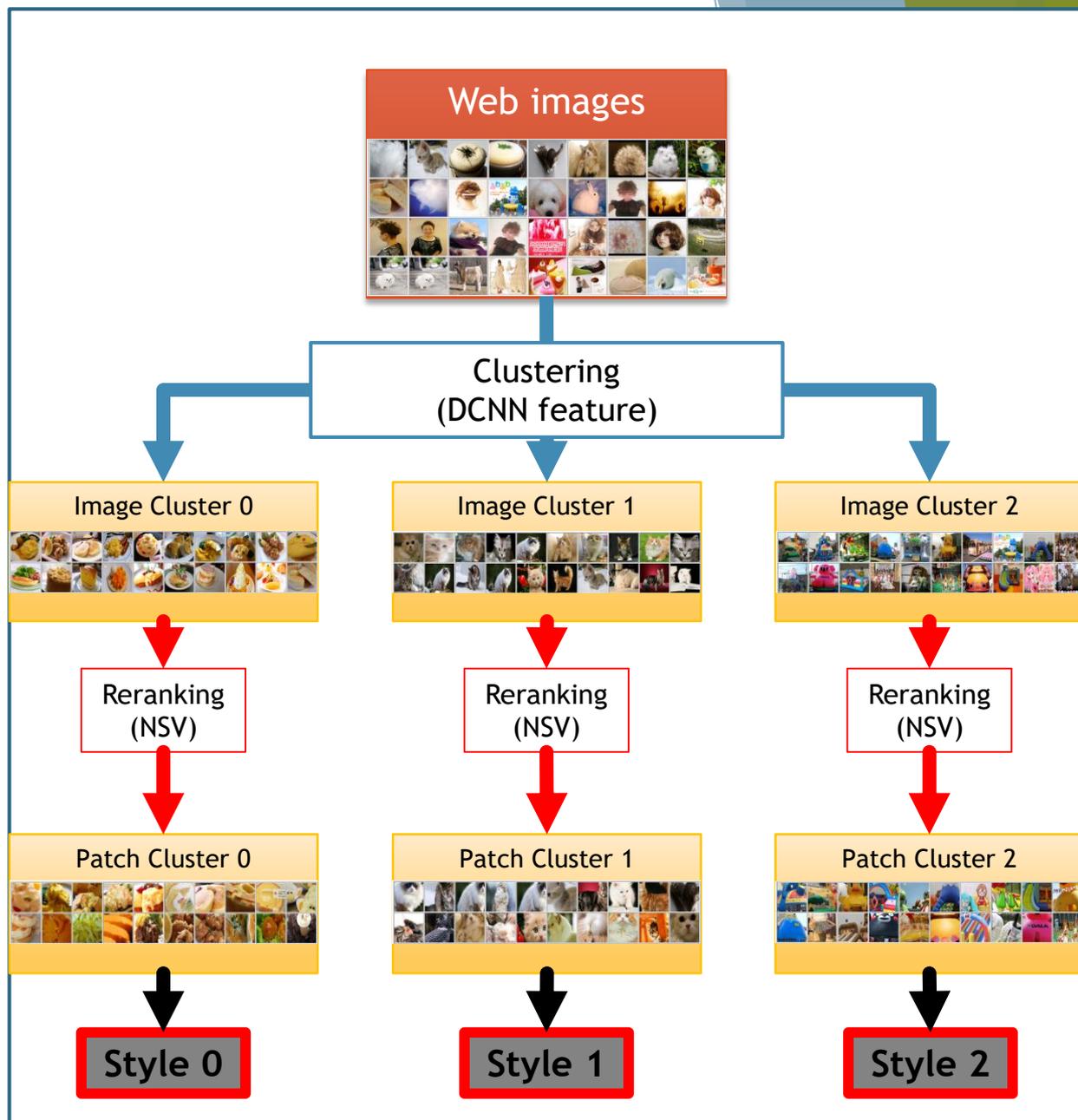
- ▶ CNNを用いた画像内物体の任意質感生成
- ▶ 質感単語からスタイル表現を自動生成



# 単語による画像 収集とクラスタ リング

Web画像をスタイルの  
類似したものの同士のク  
ラスタに分化

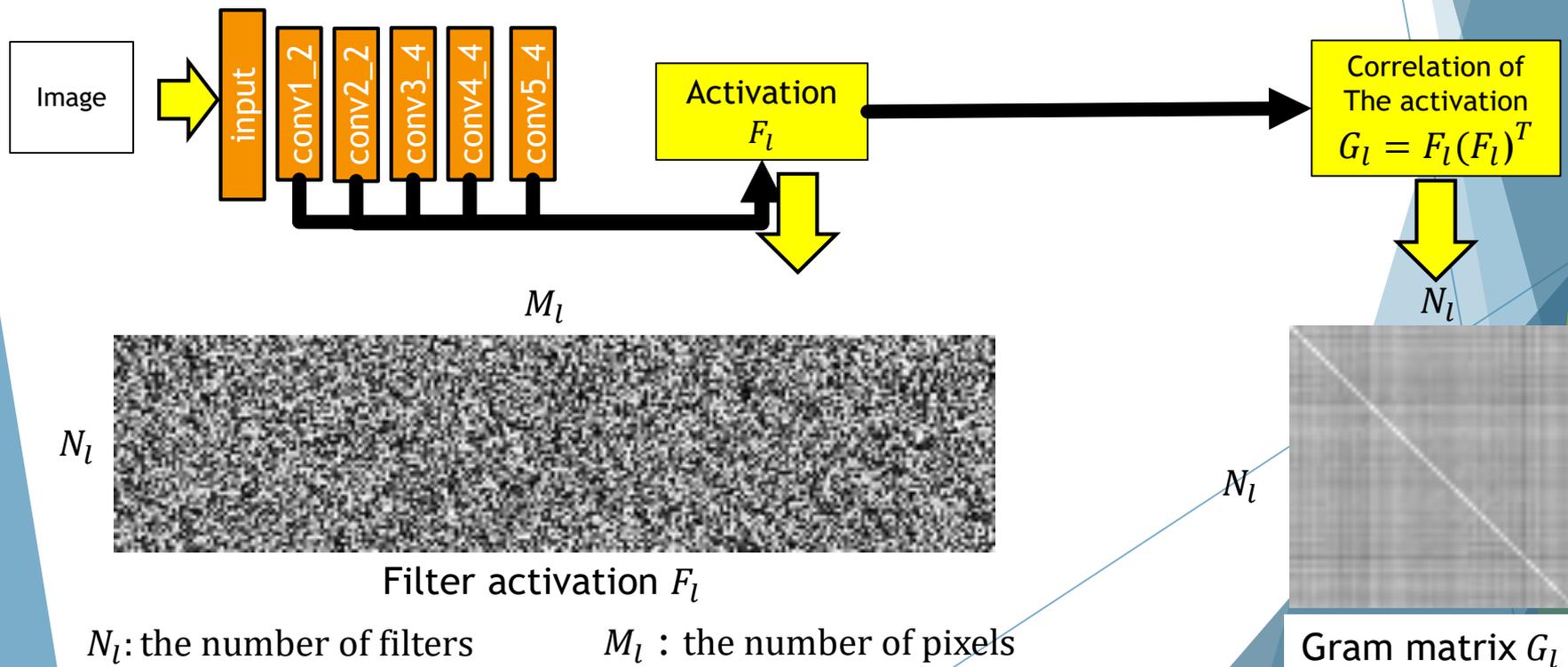
- DCNN特徴(VGG19  
fc6の出力)とNeural  
Style Vectorを使用



# Neural Style Vector

CNN-based style vector for style image retrieval. [S. Matsuo, ACM ICMR 2016.]

- ▶ スタイル変換アルゴリズムにおけるスタイル表現をベクトル化
- ▶ 絵画画像の検索・クラス分類に利用



# Neural Style Vector

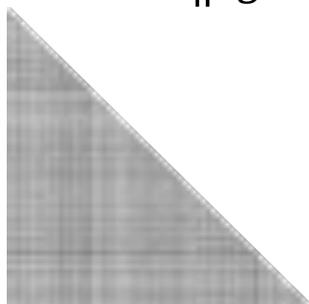
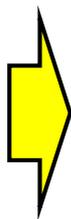
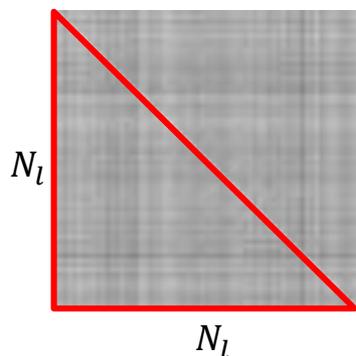
$G^l$ の対象要素と対角要素で定義

$$V^l = [G_{1,1}^l, G_{2,1}^l, G_{2,2}^l, \dots, G_{N_l,1}^l, G_{N_l,2}^l, \dots, G_{N_l,N_l}^l]$$

$$|V^l| = (\text{hurf elements}) + (\text{diagnal elements}) = N_l * (N_l + 1)/2$$

符号付平方根とL2正規化をして使用

$$V^{lsgnsqrt} = \frac{\text{sgn}(V^l)\sqrt{|V^l|}}{\|\text{sgn}(V^l)\sqrt{|V^l|}\|}$$



**Style vector  $V^l$**

Ex, at conv5\_1

$$N_{conv5_1} = 512,$$

$$|V^{conv5_1}| = 131,328$$

# 実験

## 実験データ

- ▶ コンテンツ画像：3枚
- ▶ 質感単語：オノマトペ 23 単語
- ▶ 画像収集：Bing API
- ▶ 使用クラスタ数：3

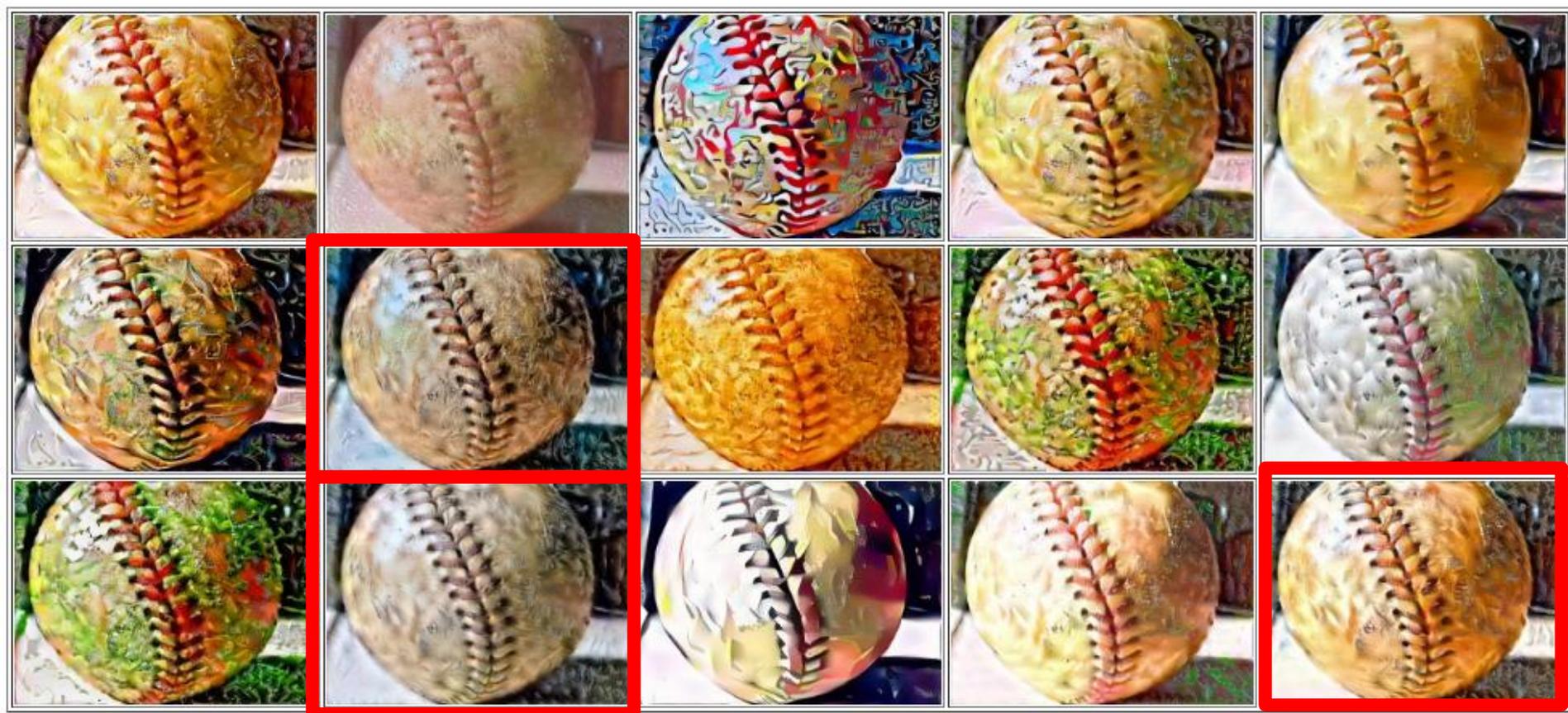
## 評価方法

- ▶ ユーザー評価
- ▶ 問題数はコンテンツ数 3単語数 23 語の 69 問
- ▶ 回答者数は 9 人

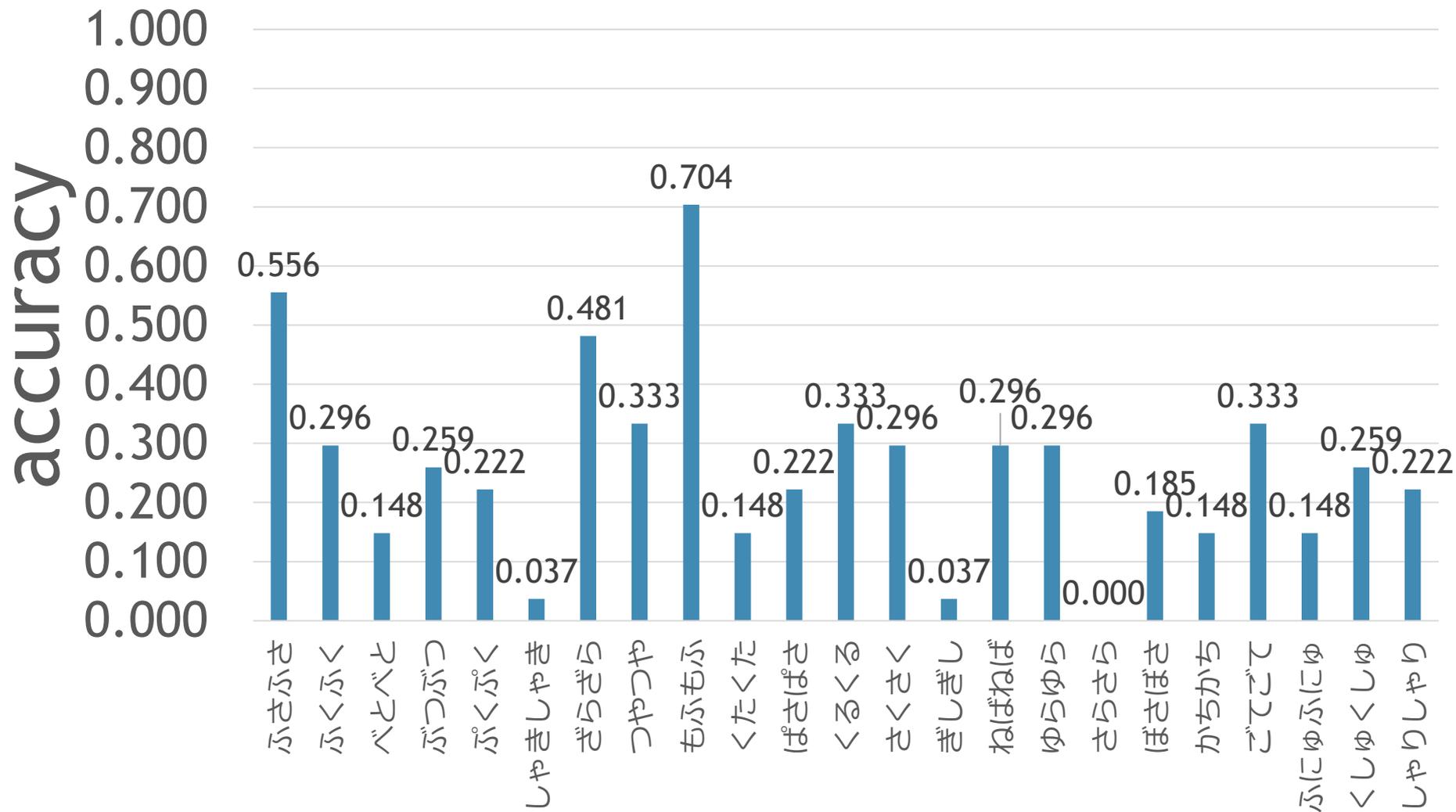
# ユーザー評価

▶ 赤枠が正解

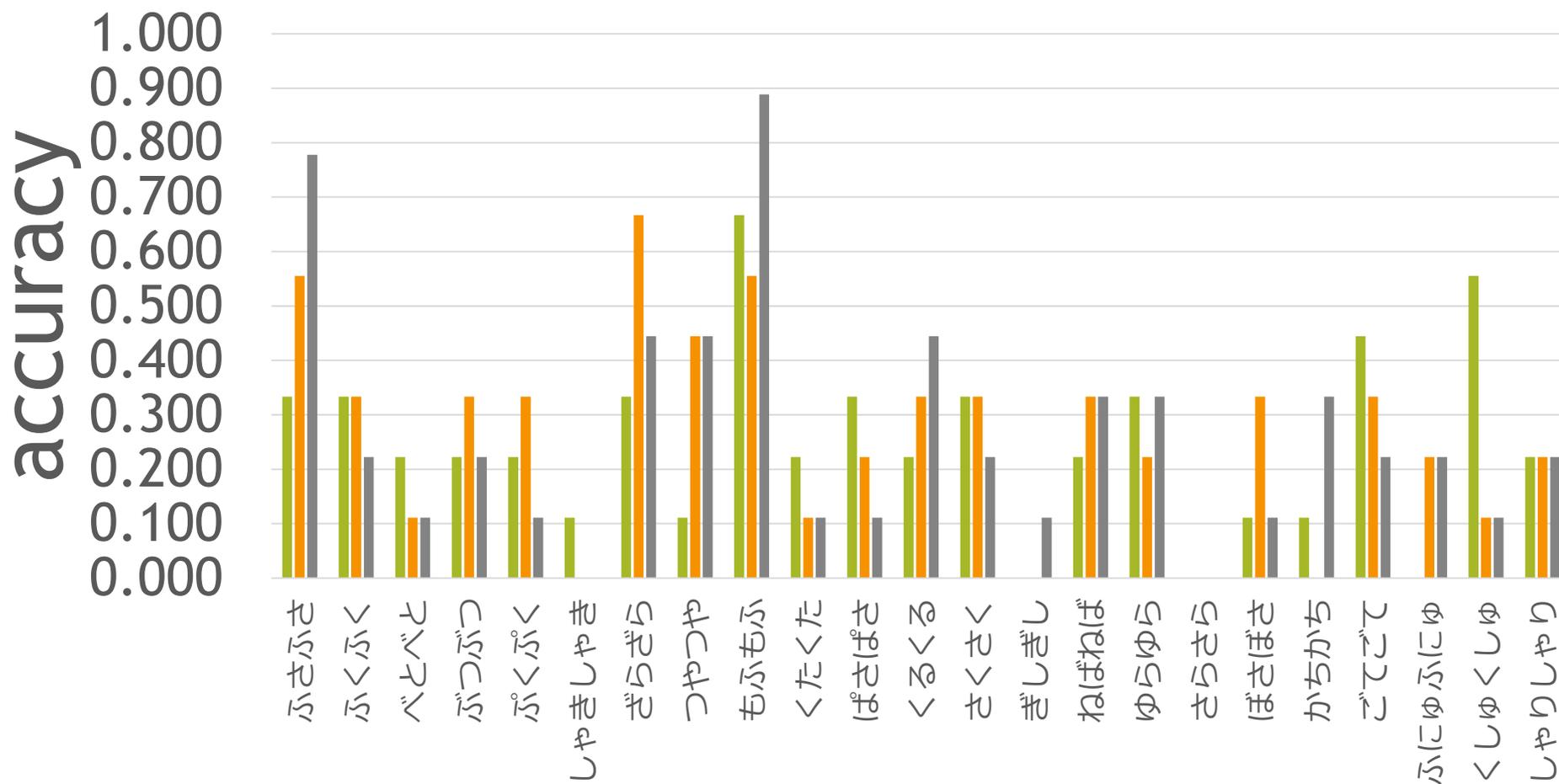
最も「ふさふさ」な画像を選んでください。0/72



# 実験結果 (総合)



# 実験結果 (コンテンツ別)



content0



content1



content2

# 結果

- ▶ 「ふさふさ」, 「ざらざら」, 「もふもふ」の正解率が高くなった。
- ▶ 他の単語では正解率は低くなった。
  1. Web画像が不適切
  2. オノマトペから画像を連想しにくい
- ▶ 各単語の正解率はコンテンツによって大きく異なっていた。
  - ▶ スタイル変換のクオリティはコンテンツ画像とスタイルの相性に強く依存していると考えられる。



Cluster 0



Cluster 1



Cluster 2



Vote: 2



Vote: 3



Vote: 5



Vote: 1



Vote: 0



Vote: 2



Vote: 0

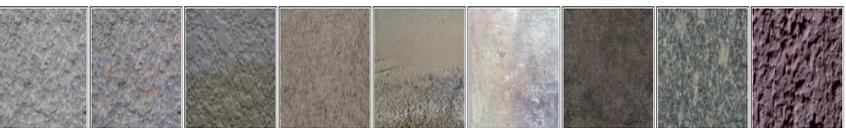


Vote: 2



Vote: 0

# 「ふさふさ」 クラスタ



Cluster 0



Cluster 1



Cluster 2



Vote:3



Vote:6



Vote:4



Vote:0



Vote:0



Vote:0



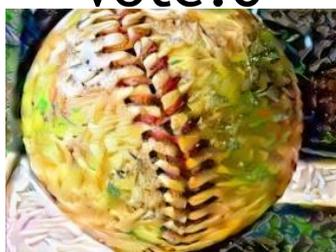
Vote:0



Vote:0



Vote:0



「ざらざら」 クラスタ



Cluster 0



Vote:0



Vote:0



Vote:0



Vote:0



Vote:0



Vote:0



Cluster 1



Vote:1



Vote:1



Vote:1



Cluster 2



Vote:4



Vote:4



Vote:7

# 「もふもふ」 クラスタ

# 考察

## コンテンツ画像と スタイルクラスタの相性

Style

Content



ふさふさ cluster0



ざらざら cluster0



もふもふ cluster2



Vote:2

Vote:3

Vote:5



Vote:3

Vote:6

Vote:4



Vote:4

Vote:4

Vote:7

# 今後の課題

- ▶ 質感単語とコンテンツの関連度と印象操作への影響
- ▶ コンテンツの持つカラー傾向や局所的構造を考慮したスタイルクラスタの構築
  - ▶ コンテンツ画像をもとにしたスタイル画像のカラー変換
  - ▶ 前景領域の利用