

深層学習による美的評価エンジンの開発と 構図推薦カメラへの実装

井上 義隆^{1,a)} 松村 択磨^{2,b)} 深澤 佑介^{1,c)} 山田 和宏^{1,d)}

概要: 本研究では深層学習による美的評価エンジンの開発と構図推薦カメラの実装を行なった。美的評価エンジンは画像の3値分類モデルとして構築した。学習済みモデルを用いた美的評価エンジンのファインチューニングによってロバスト性を獲得し分類精度を改善した。また、深層学習の推論をリアルタイムに動作させ、撮影すべき構図を撮影者に推薦するデジタルカメラ型デバイスを実装した。構図推薦カメラは初心者の撮影技術をサポートするだけでなく、熟練者に対してもシャッターチャンスの気づきや構図を追求する機会を与えることができる。

Development of Aesthetics Evaluation Engine based on Deep Learning and Implementation into Composition Recommendation Camera

INOUE YOSHITAKA^{1,a)} MATSUMURA TAKUMA^{2,b)} FUKAZAWA YUSUKE^{1,c)} YAMADA KAZUHIRO^{1,d)}

1. はじめに

近年、デジタルカメラに内蔵されている受光センサや手ぶれ補正技術の性能向上により、照明が暗い環境下であっても、ノイズを少なく抑えて、手ぶれによる画像の品質劣化を発生させることなく撮影することが可能になった。また、顔や瞳の検出とオートフォーカスによって、人物の表情を確実に捉えることが可能になった。さらに1秒あたりの撮影可能枚数の増加は動物、スポーツ等の動体を被写体とした場合にシャッターチャンスを確実に抑えることが可能になった。

デジタルカメラの性能向上により、きれいな写真を撮影することが可能になったものの、必ずしもそれだけで本質的に「美しい」と感じる写真を撮影できるとは限らない。

人間が写真を「美しい」と判断する際には、正確な露光やピントという光学的観点だけではなく、写真の構図や内容に基づく情緒的観点も影響する。後者は人物の魅力的な表情、大自然の中の非日常的な絶景、構図の幾何学的な精密さ、二度と訪れない決定的瞬間等が挙げられる。色彩を考慮しても、街中の色鮮やかな看板広告よりは、花畑や夕焼けの方が一般的に好まれるかもしれない。モノクロ写真では、コントラストが強く幾何学的な写真も、ノスタルジックな淡い写真も等しく好まれるかもしれない。

従来の画像処理技術では情緒的観点の特徴量化することは困難であるため、大量の美的評価済み画像から深層学習によって光学的観点と情緒的観点とを同時に抽出する手法が研究されている。しかしながら、既存手法は美的評価を2値分類問題として取り扱っており、多段階の評価を行っていない。2値分類問題の場合、美的品質のわずかな改善を見逃してしまう恐れがあり、撮影技術向上を目的とした活用には適していない。また、既存手法はトレーニングとテストとでデータを分割しているものの、与えられたデータセットの中での分類精度の最適化に留まっている。製品として提供する場合はユーザが撮影する写真を事前に想定できないため、性質の異なる入力画像に対するより強

¹ 株式会社 NTT ドコモ
KOKUSAI AKASAKA building, 4-5, Akasaka 2-chome,
Minato-ku, Tokyo 107-0052, Japan

² ドコモ・テクノロジー株式会社
DOCOMO R&D center, 3-5, Hikarino-oka, Yokosuka-shi,
Kanagawa 239-8536, Japan

a) yoshitaka.inoue.ye@nttdocomo.com

b) matsumuratak@nttdocomo.com

c) fukazawayu@nttdocomo.com

d) yamadakazu@nttdocomo.com



図 1: 構図推薦カメラ本体. (上) 正面, (下) 背面, (右) 内部
Fig. 1 Composition recommendation camera. (Top) Front, (Bottom) Back, (Right) Inside.

いロバスト性が必要となる. またこれらの手法をデジタルカメラ型のデバイスに実装し, 撮影者に対してリアルタイムに美的評価と構図推薦を行う製品は存在しない.

本研究では, 美的評価済みの AVA データセット [1] を用いて深層学習による美的評価エンジンを開発し, 写真の美しさを 3 段階で判別することを可能にした. 撮影済み写真に対してだけではなく, 撮影現場でリアルタイムに美的評価を実行できるように, エンジンを軽量の GPU マシン上に実装し, レンズとセンサと組み合わせて, 手持ちで撮影できる構図推薦カメラを開発した (図 1). 構図推薦カメラはライブビュー画像を逐次評価するだけでなく, より高い評価値が得られる構図を撮影者に推薦することができる.

本研究の貢献ポイントは以下の通りである. まず, Accuracy 70.0% の 3 段階の美的評価を可能にした. 次に, セマンティクス分類で学習済みのモデルをもとに美的評価エンジンをファインチューニングすることで, 入力画像に対するロバスト性を獲得するとともに分類精度を改善した. また, 撮影機能と美的評価をデバイスローカルで完結させ, リアルタイムな動作を可能にした. さらに, 撮影初心者が犯しがちな構図ミス指摘したり, 熟練者に対しても気づきや構図を追求する機会を与えるような構図推薦機能を実装した.

以下, 2 章では関連研究について述べ, 3 章で構図推薦カメラに求められる要求条件とアプローチについて述べる. 4 章では美的評価エンジンの構築, 5 章で構図推薦カメラの実装について述べる. 6 章で美的評価エンジンと構図推薦カメラの評価について述べる. 最後に 7 章でまとめる.

2. 関連研究

2.1 特徴量設計によるアプローチ

写真を評価する際に構図は重要な要素である [2], [3], [4]. 代表的な構図ルールとして日の丸構図, 三分割構図, 黄金分割構図, 対角構図等が挙げられる. このような構図ルールに従った写真は美しさや安定感を与えられると言われている.

構図や色を特徴量として抽出し, ルールベースで写真を美的評価する手法が研究されている. 家田らは入力された写真を解析し, 色による顔の検出, 色彩的に際立つ領域, 三角形, 水平線, 対角線, 遠近法消失点を考慮して, 構図ルールに近づくようにトリミング領域を推薦する手法を提案している [5]. 志津野らの手法は入力された写真から SURF 特徴量を抽出し, あらかじめ用意した構図ルールのテンプレートと比較して, 従うべき適切な構図ルールを撮影者へ推薦する [6]. 撮影者が当該構図ルールに合わせるように撮影することで良い写真が得られる. Bhattacharya らは構図ルールに基づいた特徴量を抽出し, 別途用意した 632 枚の美的評価済み写真を用いて, Support Vector Regression によって美的評価値を推定するモデルを構築している [7].

情緒的観点において「美しい」写真であるための要素は, 構図ルールだけではなく写真の内容も含まれる. 画像解析によって構図ルールのような幾何学的特徴量を抽出することは可能であっても, 写真の内容を特徴量として抽出することは困難である.

2.2 Deep Learning によるアプローチ

AVA データセットは 25 万枚の写真を含み, 1 枚当たり 78 人から 549 人の評価者によって 10 段階の主観評価が記録されている [1]. AVA データセットを用いた深層学習による美的評価手法が研究されている.

Lu らの手法は, Convolutional Neural Network (以下, CNN) の入力層が規定する解像度になるように入力画像を変換し, 画像の全体と細部とのそれぞれを分割して同時にネットワークへ入力する [8]. 従来の CNN は入力層が規定するサイズに入力画像をリサイズしなければならないため, 画像が歪むという問題がある. しかしながら, 入力画像の解像度やアスペクト比は, 撮影条件やカメラの設定, 撮影後の編集に依存するためあらかじめ想定することができない. Lu らは次の手法でこの問題を解決している. 画像の全体の成分として, 入力画像をクロップ (入力画像の短辺の長さで中央を正方形をトリミング), ワープ (長辺のみを縮小), パディング (長辺を合わせて, 短辺方向に生じる隙間を黒で埋める) し, さらにそれぞれを正方形の入力層に合わせてリサイズする. これは構図やグラデーションを情報として含む. また, 画像の細部の成分として, 画像の部分領域をランダムに選択し, 入力層に合わせてトリミングする. これは画像のテクスチャ情報を残している. それぞれを個別のネットワークに入力し, 最終レイヤーで統合する. AVA データセットに対する 2 値分類 (高スコア, 低スコア) で精度 74.46% を達成している.

Kao らは画像のセマンティクスを考慮して CNN で美的評価を行っている [9]. セマンティクスとは画像の内容が人物なのか, 建築物なのか, 食べ物なのか, といった情報である. セマンティクスを考慮することで 2 値分類精度を

79.08%まで向上させている。

Luらの手法 [8] と Kaoらの手法 [9] は美的評価を2値で分類しているが、撮影技術向上の目的でより有益な多段階評価を可能にしていない。また、AVA データセットはアマチュア写真家がオンラインコミュニティに投稿した25万枚の画像から構成される。すなわち、投稿者によってあらかじめ選別された画像であるため、カメラの通常のユースケースで入力される画像と比較して、内容に偏りがある恐れがある。本研究が提案する構図推薦カメラを製品として提供するには入力画像に対するロバスト性が重要となる。

3. 問題設定

3.1 要求条件

本研究は美的評価の精度をある程度に保ちつつ、撮影作業をリアルタイムにかつインタラクティブにサポートするデバイスの開発を目的とする。これにより初心者向けには撮影技術の向上を、熟練者向けにはシャッターチャンスの気づきや良い構図を追求する機会を与える。具体的に本研究で実現する美的評価エンジンおよび構図推薦カメラの要求条件は以下の通りとする。ただし、以下ではリアルタイム性の要求条件を定義する際に図2を参照するが、図2の詳細は5.5節および5.4節で述べる。

多段階評価 3段階以上の美的評価を可能にする。撮影技術の向上のためには、写真のわずかな改善に対しても鋭敏に美的評価を反映させる必要がある。

ロバスト性 製品として提供する場合はユーザが入力する画像は想定できないため入力画像に対するロバスト性が必要となる。

リアルタイム性 画像の取得と表示の処理は25.0fps以上とする。図2の(a) → (b) → (c)のループに該当する。撮影者がスムーズに構図を探索し、シャッターチャンスを逃さないよう、実風景と表示内容にラグが存在してはならない。また、フレームレートは標準的なBlu-ray Disc映像と同等であることを基準とする。推論処理は2.5fps以上を目標とする。図2の(e) → (f) → (g) → (h)のループに該当する。実環境が1秒の間に大きくは変化しないと仮定している。

構図推薦 実風景のうち、撮影画角の周辺の画像を含めて美的評価の高い写真が得られる構図を探索し、撮影者に推薦する。撮影技術向上のためには撮影者が気づかない構図を撮影者に提示する必要がある。

3.2 アプローチ

前節の要求条件を次のアプローチによって解決する。まず、AVA データセットの各画像のスコアをもとに3つのラベルに分類し、3値分類問題としてモデルを構築する。次に、異なるデータセットで学習済みのモデルをもとに美的評価エンジンをファインチューニングすることでロバス

ト性を獲得する。また、リアルタイム性については以下の3つの要素から実現する。(1) 深層学習の際にシンプルなネットワーク構造を用いる。(2) 通信遅延を発生させないためにデバイスローカルで処理を完結させる。(3) 表示と推論を並列プロセス処理させることでスムーズな操作性を実現する。さらに、構図推薦については、ライブビュー画像を重畳するように9枚の画像に分割して、9枚の画像に対して同時に美的評価を行い、評価結果をもとに最適構図で撮影するための移動方向を撮影者に提示することで解決する。

上記それぞれのアプローチについては、3値分類モデルを4.1節で述べ、ファインチューニングを4.3節で述べ、リアルタイム性のための(1)ネットワーク構造を4.2節、(2)デバイスを5.2節、(3)並列プロセス処理を5.4節でそれぞれ述べ、構図推薦アルゴリズムを5.5節で述べる。

4. 美的評価エンジンの構築

4.1 3値分類モデルの構築

多段階評価の要求条件を満たすために下記の通りデータセットを整理して3値分類問題として定義する。

AVA データセット [1] は、画像のインデックスを $i \in I$ 、画像を X_i として、画像 X_i それぞれについてスコア $\{s \in \mathbb{N} \mid 1 \leq s \leq 10\}$ を付与した評価者の人数分布が $n_{i,s}$ として記録されている。画像 X_i ごとに評価者人数 $N_i = \sum_{s=1}^{10} n_{i,s}$ は異なり、AVA データセットにおいては最小で78人、最大で549人であった。画像 X_i ごとに評価者で平均したスコア $S_i = \frac{1}{N_i} \sum_{s=1}^{10} n_{i,s} s$ を算出する。画像で平均したスコア $\mu = \frac{1}{|I|} \sum_i S_i$ と、標準偏差 $\sigma = \sqrt{\frac{1}{|I|} \sum_i (S_i - \mu)^2}$ とを基準として、ラベル $Y_i \in \{\text{High, Middle, Low}\}$ を付与する。具体的にはラベル High は $\mu + \sigma \leq S_i$ 、Middle は $\mu - \sigma \leq S_i < \mu + \sigma$ 、Low は $S_i < \mu - \sigma$ を満たす画像 X_i に付与する。

分類された画像の枚数はラベルごとに異なる。いずれのラベルも画像の枚数が同数程度になるようにランダムにサンプリングし、サンプリングされた画像を90%と10%に分割して、それぞれをトレーニングデータ、テストデータとする。ただし、汎化性能を向上させるために、トレーニングデータについては画像 X_i を水平方向に反転した画像も追加する。被写体が左右非対称であることを事前知識として持っている場合(看板の文字、衣服のボタン等)、見た目違和感を受けるかもが、そのような写真の数はごく一部であって、動植物、人物、建築物、都市や自然の風景については左右反転しても気づかず、美的品質に影響しないと考えられる。なお、テストデータには水平方向に反転した画像は含まれない。ラベルごとの画像枚数を表1に示す。左列から、元の画像枚数、トレーニングデータの画像枚数、テストデータの画像枚数である。

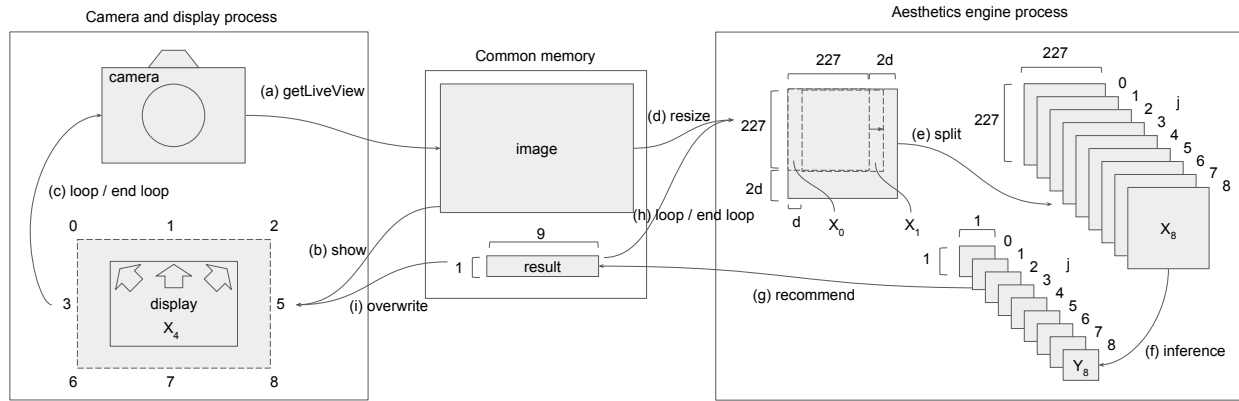


図 2: 構図推薦カメラと美的評価エンジンのアルゴリズム

Fig. 2 Algorithm of composition recommendation camera and aesthetics evaluation engine.

表 1: データセットの 3 値分類定義

Table 1 Definition of dataset classified by three labels.

Label	#original images	#training images	#test images
High	39,580	68,959	3,888
Middle	177,643	68,951	3,784
Low	38,307	68,979	3,839

4.2 CNN

リアルタイム性の要求条件を満たすための要素 (1) として、複雑なネットワーク構造は用いずに ImageNet[10] のネットワーク構造を採用する。

前処理にかかる計算量を削減するために、画像のワープのみ実行して入力層に合わせる。長辺方向を縮小して、縦 227, 横 227, 3 チャンネルのサイズに統一する。

また、本研究では ImageNet のネットワーク構造の出力層の次元を 3 に変更する。4.1 節で定義したラベルの 3 値分類を学習する。出力される値は画像 X_i ごとの、3 種類のラベルそれぞれへの所属確率 $P(Y_i | X_i)$ となる。また、1 枚の画像について 3 カテゴリへの所属確率の和は $\sum_{Y_i} P(Y_i | X_i) = 1$ となる。

4.3 ファインチューニング

ロバスト性の要求条件を満たすために、ファインチューニングによって解決する。

1400 万枚の画像を含む ILSVRC2012 データセットを用いて、セマンティクス観点での分類を学習したリファレンスモデル (caffe_reference_imagenet_model) が公開されている [11]。ILSVRC2012 は画像の枚数が AVA データセットよりはるかに多いこと、かつアマチュア写真家による投稿画像以外の幅広い内容の画像をカバーしていることから、リファレンスモデルを初期値として美的評価を学習、すなわちファインチューニング (以下、FT) することで、ロバ

スト性を獲得できると考えられる。

4.4 モデル構築環境

Amazon Web Service の EC2 インスタンスを用いる。インスタンスタイプは g2.2xlarge[12] とし、Ubuntu にインストールした Caffe を用いる。

5. 構図推薦カメラの実装

5.1 概要

4 章で構築した美的評価エンジンを構図推薦カメラとしてデバイスに実装する。ハードウェア構成、並列プロセス処理、構図推薦処理を以下で述べる。

5.2 ハードウェア構成

リアルタイム性の要求条件を満たす要素 (2) として、美的評価エンジンを NVIDIA Jetson TX1[13] (以下、TX1) へ実装し、ローカルで推論処理を行う。

TX1 は NVIDIA GPU Tegra X1 を搭載した小型軽量のコンピュータである。CPU, GPU, メモリ, ストレージを含む TX1 モジュールと、Wi-Fi, HDMI 端子, GPIO, 電源等を含むキャリアボードとで構成されている。TX1 モジュールを別売りの小型のキャリアボード Orbitty Carrier for NVIDIA Jetson TX1[14] (以下、Orbitty Carrier) に載せ替えた。Orbitty Carrier は縦 5cm, 横 8.5cm である。

Orbitty Carrier と LiPo バッテリ, 小型の液晶ディスプレイを接続した。アクリル板を加工してカメラの筐体を作成し、TX1, LIPO バッテリを筐体内に収納し、ディスプレイは筐体背面側に固定した。図 1 右図は液晶ディスプレイを取り外した状態の筐体内部である。筐体内部の左側が TX1 モジュールと Orbitty Carrier である。

レンズユニットは SONY ICLE-QX1[15] (以下、QX1) を用いた。QX1 はセンサとシャッターボタンのみ備えており、ライブビューを確認するためのディスプレイを搭載

しておらず、設定変更操作に必要なダイヤルやボタンも備えていない。本研究では QX1 をレンズユニットとして採用し、QX1 と TX1 を接続し、TX1 と液晶ディスプレイを接続した。なお、レンズは SONY SEL35F18 (以下、レンズ) を用いた。

QX1 は筐体外側に固定した。操作性を考慮して、QX1 に備わっているシャッターボタンとは別に新たなシャッターボタンを作成して筐体外側に配置した。またボタン操作検知用マイコンボードを筐体内に配置し、GPIO で Orbitty Carrier と接続した。図 1 上下図のように、QX1 およびレンズが筐体に固定されている。また、図 1 下図の右下に見えるシルバーのボタンがシャッターボタンである。

以上のハードウェア構成で、筐体の幅 19.0cm、高さ 9.5cm、厚さ 6.3cm、ディスプレイからレンズの先端までの長さ 21.0cm、全体の重量は 1082g となった。デジタルカメラのハイエンド機のと比較しても重量差はあまりない。SONY α 7 II のズームレンズキットはバッテリーとレンズ込みで 894g、Canon EOS 5D Mark IV のレンズキットはバッテリーとレンズ込みで 1490g である。

5.3 ソフトウェア環境

TX1 の Ubuntu 上に、QX1 との通信、CNN による推論、表示等の処理を Python で実装した。QX1 は SONY Camera Remote API に対応している [16]。QX1 内部に Wi-Fi アクセスポイントを立ち上げ、TX1 から Wi-Fi 接続することで、QX1 で取得したライブビューを取得したり、QX1 に対してシャッター命令を行うことができる。このような処理は HTTP 通信で行われる。TX1 上の Python から QX1 と通信できるようにライブラリを構築した。また、画像処理と表示は OpenCV を、推論は Caffe を用いた。

5.4 並列プロセス処理

リアルタイム性の要求条件を満たす要素 (3) として、並列プロセス処理と共有メモリによるアプローチをとる。

撮影現場での活用のためには、カメラの指示に従って撮影者がカメラを動かし、カメラに新しい画像が入力されてカメラが撮影者への指示を変更する、というような撮影者とカメラとの間のインタラクティブなやりとりを遅延なくスムーズに実現しなければならない。

そのため本研究では、図 2 のように画像の取得および表示をコントロールする表示系プロセス (Camera and display process) と、美的評価を逐次行う推論系プロセス (Aesthetics engine process) とを分割し、共有メモリ空間 (Common memory) を設置する。もし、2つのプロセスを同一のプロセスで実行すると、推論の際に負荷がかかり、表示画像が固まるようなラグが生じてしまう。例えば、構図の誘導指示に従ってカメラをスライドさせる際に、ラグが生じるたびにディスプレイ上のライブビューが一時静止してしまい、

滑らかに動作しない。プロセスを分割することでスムーズな動作を可能にする。

表示系プロセスは、図 2 (a) のように取得したライブビュー画像を共有メモリに格納し、図 2 (b) のようにディスプレイにライブビュー画像を表示する。このとき、図 2 (i) のように共有メモリに推論結果が格納されていればそれをライブビュー画像に重畳表示する。図 2 (c) のように表示後には再度ライブビュー取得へ戻る。

一方、推論系プロセスは、図 2 (d) のように共有メモリに画像が格納されていれば、その画像に対して美的評価を行う。また、図 2 (g) のように推論結果を共有メモリに格納する。格納後は再度共有メモリ上の画像確認に戻る。

5.5 構図推薦処理

構図推薦の要求条件を満たすために、下記の構図推薦処理を実装する。

図 2 (d) に示すように、QX1 から取得した画像を $227+2d$ 四方に縮小させた後に、図 2 (e) に示すように、幅 d でスライドしながら重複するように画像を 9 分割する。227 とは 4.2 説で定義した CNN の入力画像サイズである。9 分割した画像のインデックスを $\{j \in \mathbb{N} \mid 0 \leq j \leq 8\}$ 、各画像を X_j とする。ただし、左上から水平方向優先で走査順にインデックスを割り当てるものとし、例えば左上端を $j = 0$ 、右上端を $j = 2$ 、中央を $j = 4$ とする。

次に、図 2 (f) に示すように、美的評価エンジンに 9 枚同時に入力し、ラベル 3 値への所属確率 $P(Y_j \mid X_j)$ を出力し、式 1 で総合スコアを算出する。

$$\text{TotalScore}_j = \sum_{Y_j} a(Y_j)P(Y_j \mid X_j) \quad (1)$$

ただし、 $(a(\text{High}), a(\text{Middle}), a(\text{Low})) = (1.0, 0.5, 0.0)$ とする。

図 2 (g) に示すように、評価結果に基づく推薦方向を result に格納する。具体的には、中央の画像 X_4 の総合スコア TotalScore_4 があらかじめ設定した閾値 th_{OK} 以上である場合、もしくは中央の画像 X_4 の総合スコア TotalScore_4 が、9 枚の中で最大である場合は、推薦方向 result = [4] とする。この場合には図 2 (i) で現在の構図のまま撮影を指示する。より高い総合スコアを得られる方向が X_4 以外に 1 個以上存在する場合は、それらの方向をすべて推薦方向 result とし、図 2 (i) でライブビュー画像に矢印で重畳表示して最適構図への移動方向を推薦する。撮影者はこの矢印の方向へカメラの向きを移動させることでより美しい写真が得られる領域を探索することができる。

ただし、本来は QX1 の受光センサーで広範囲の画像を取得しているが、撮影者にとっては中央部分 X_4 のみが撮影対象となる。撮影者には X_4 だけを表示して、 X_4 以外の領域を非表示にしてもよい。

Algorithm 1 CameraAndDisplayProcess

```

1: global image ← null
2: global result ← null
3: camera.initialize()
4: aestheticsEngineProcess.start()
5: loop
6:   image ← camera.getLiveView()
7:   display.show(image)
8:   if result ≠ null then
9:     display.override(result)
10:  end if
11:  if camera.shutterButton.onPressed() then
12:    camera.takePictureAndSave()
13:  end if
14: end loop

```

Algorithm 2 AestheticsEngineProcess

```

1: cnn ← loadModel()
2: loop
3:   if image ≠ null then
4:     image ← resize(image)
5:     X[0, ..., 8] ← split(image)
6:     Y[0, ..., 8] ← cnn.inference(X)
7:     result ← recommend(Y)
8:   end if
9: end loop

```

5.6 構図推薦カメラのアルゴリズム

5.4節と5.5節の処理の全体をアルゴリズム1, 2に示す。表示系プロセス (CameraAndDisplayProcess) 内部でライブビュー取得画像 image と推薦結果 result を、推論系プロセス (AestheticsEngineProcess) との共有メモリ上の変数として定義する。また表示系プロセスは推論系プロセスを起動 (start) する。表示系プロセスのループでは、image の取得 (getLiveView), 表示 (show) を行い、result が存在する場合には推薦結果を重畳表示する (overwrite)。また、シャッターボタンが押された場合は、撮影して記録保存する (takePictureAndSave)。

推論系プロセスは美的評価を行う。推論系プロセスのループでは、新しい image が存在する場合は画像を重畳分割 (split) 後に推論 (inference) を行い、推論結果に基づいて計算 (recommend) した推薦結果を result に格納する。

以上の2つのプロセスが image と result を逐次更新する。

6. 評価

6.1 多段階評価およびロバスト性の評価

本節では多段階評価の精度と FT による改善を確認する。ラベル $Y_i \in \{High, Middle, Low\}$ への所属確率 $P(Y_i | X_i)$ のうち、最大所属確率 $\hat{Y}_i = \arg \max_{Y_i} P(Y_i | X_i)$ をとるラベル \hat{Y}_i を推定ラベルとする。テストデータ画像 X_i の正解ラベル Y_i と、推定ラベルの \hat{Y}_i の組合せ (Y_i, \hat{Y}_i) ごとに枚数を集計した結果を図3に示す。左が CNN のみ、右が

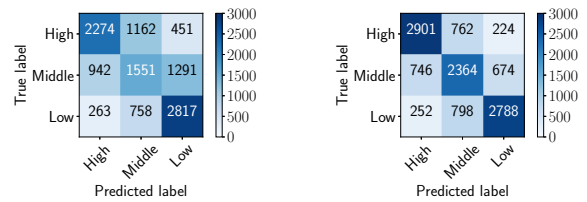


図3: 混同行列。(左) CNN, (右) CNN+FT
Fig. 3 Confusion matrix. (Left) CNN, (Right) CNN+FT.

表2: 精度評価

Table 2 Evaluation by accuracy/precision/recall.

		CNN	CNN+FT
Accuracy		57.7%	70.0%
Precision	High	65.4%	74.4%
	Middle	44.7%	60.2%
	Low	61.8%	75.6%
Recall	High	58.5%	74.6%
	Middle	41.0%	62.5%
	Low	73.4%	72.6%

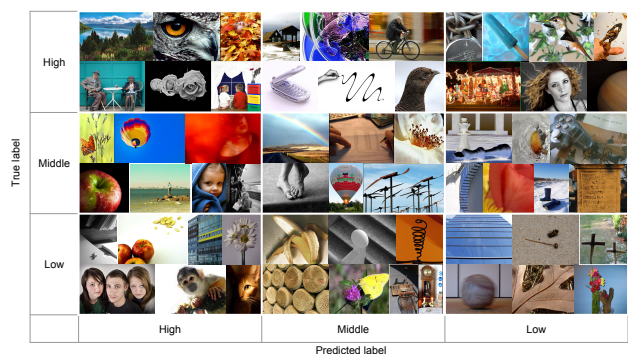


図4: 分類された画像の一例

Fig. 4 Examples of classified images.

CNN+FT の混同行列である。対角成分の値が大きいほど分類精度が高い。CNN+FT は、美的評価をより正確に推論できていることが確認できた。また、CNN に比較して CNN+FT は特に (High, High) および (Middle, Middle) を大きく改善した。

図3をもとに分類精度を算出した結果を表2に示す。Accuracy は CNN の場合 57.7% だったが、CNN+FT では 70.0% まで改善した。Recall の Low を除き、CNN+FT によって各指標は大幅に改善している。本節では FT によってロバスト性を獲得して精度が改善し、Accuracy 70.0% で3値分類可能なことを確認した。

CNN+FT で分類された画像の例を図4示す。推定ラベルが High である写真は、色合いが鮮やかで、コントラストが強く、被写体の内容がはっきりしているもの、構図が練られている写真が多い印象を受けた。推定ラベルが Low である写真はブレやノイズを含み、被写体の内容が不鮮明である写真が多い印象を受けた。

6.2 リアルタイム性の評価

本節では構図推薦処理のリアルタイム性を確認する。

QX1が取得するライブビューの解像度は高さ424、幅640であるが、重複する幅を $d=90$ として、407四方に縮小した。また、中央画像での撮影を指示する閾値を $th_{OK}=0.9$ とした。並列プロセス処理により、表示系プロセスは常に約25fpsで処理され、推論系プロセスの負荷による表示ラグが発生することはなかった。推論系プロセスは最大で約5.0fps（ライブビュー画像5フレーム毎に推論）で処理可能であったが、被写体の状況が1秒間以内に大きく変化することはないと想定して、実装上は約2.5fps（10フレーム毎）に抑えて動作させた。以上の処理速度からリアルタイム性の要求条件を満たすことを確認した。

6.3 構図推薦の動作例

本節では構図推薦の挙動を確認する。

推薦表示の例として、特に期待された効果を示した例を図5に、期待以上の気づきが得られた例として図6に示す。いずれの図も構図推薦カメラのライブビュー画像である。5章で述べたように、中央の緑色の矩形内部が画像 X_4 、すなわち撮影者が意識する撮影画角を示している。 X_4 に重畳するように赤色の三角形で構図推薦方向を示し、「OK」という赤色のテキストで撮影タイミングを指示している。また、 $TotalScore_4$ が0.75以下の場合には「GOOD」、0.25未満の場合には「BAD」と緑色で表示している。

図5は期待された効果を示したポートレート撮影の例である。撮影者から見て左に男性が、右に女性が立っている。図5の左図の状況では、男性の顔の上半分と女性の頭部が画像 X_4 の外側に見切れていた。すると、右上へ構図を移動させるべきである構図推薦結果が表示された。この推薦結果に従い、撮影者が右上方向へ構図を変更したのが図5の右図の状況である。男性と女性の顔全体が X_4 に収まった瞬間に撮影指示が表示された。

図5のように、被写体が見切れている状態というのは、初心者に見られる単純な構図のミスであるが、構図推薦カメラはこれを見逃さずに指摘していることが確認できた。AVAデータセットに含まれる画像は1枚1枚異なり、同一の被写体に対して良い構図と悪い構図が含まれているわけではない。しかしながら、ラベルHighに含まれる画像の多くは良い構図で撮影されている写真が多いため、構図推薦結果もまた良い構図になったと考えられる。

図6は期待以上の気づきが得られた風景撮影の例である。奥に森林、手前に湖面がある。湖面に木々が反射している。図6の左図の状況では、 X_4 の下端近くに湖面と森林の境界線が位置していた。構図推薦は左下、下、右下への移動を示した。撮影者が下方向へ構図を調整したのが図6の右図の状況である。湖面と森林の境界線が X_4 の下端から1/3に位置した瞬間に撮影指示が表示された。

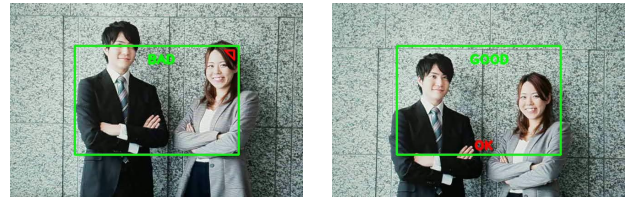


図5: ポートレート撮影での動作例。(左) 構図推薦、(右) 撮影指示

Fig. 5 Demonstration for portrait. (Left) Recommend composition, (Right) Photo opportunity.

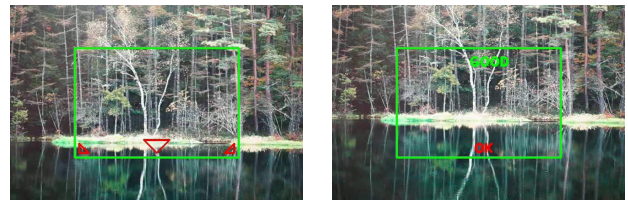


図6: 風景撮影での動作例。(左) 構図推薦、(右) 撮影指示

Fig. 6 Demonstration for landscape. (Left) Recommend composition, (Right) Photo opportunity.



図7: 3分割構図の例

Fig. 7 Rule of thirds.

図6のように、画面内を縦横3等分した線に風景の水平線や地平線を合わせたり、あるいは縦横の線の交点に被写体を配置する構図は3分割構図と呼ばれる。AVAデータセットにもともと含まれていた、撮影者自らが3分割構図に従って撮影したと思われる写真の例を図7に示す。3分割構図は熟練者が意識するような複雑な構図であるが、構図推薦カメラの推薦結果はこの構図に従っていた。本研究は明示的に3分割構図を学習したわけではないが、ラベルHighに含まれる画像の多くに3分割構図が出現していたため、3分割構図を推薦したと考えられる。

6.4 構図推薦の評価

6.3節では構図推薦の動作の一例を示したが、複数のユーザの観点で構図推薦の要求条件、すなわち、構図推薦機能によって美しい写真を失敗することなく撮影できるようになるかどうかについて評価を行う。

被験者29人に構図推薦カメラを操作してもらい、アンケートを実施した。従来のデジタルカメラに構図推薦機能を付加した場合に、追加で支払い得る金額を質問とした。

表 3: 構図推薦に対する評価アンケート結果

Table 3 Questionnaire results on composition recommendation.

Monthly fee (yen)	0	100	300	1,000	3,000
#answers	3	6	18	2	0

なお、デジタルカメラの価格帯は数万円から数十万円と幅広く、被験者の相場感覚が一定ではないことから、構図推薦機能の利用に対して月額料金を要すると仮定し、選択肢（無料、100円、300円、1000円、3000円）の中から回答してもらった。評価結果を表3に示す。選択肢の中で月額300円が最も多くの回答者を得た。全29人中、無料と回答した3人を除く26人により、構図推薦機能に一定の価値があることが認められた。構図推薦機能を備えたデジタルカメラであれば価格帯が上昇しても購入検討対象となり得ることが確認できた。

7. おわりに

本研究では、写真の美的観点から High, Middle, Low の3段階に分類する美的評価エンジンを CNN によって構築した。その際に学習済みのモデルを参照して、美的評価エンジンのファインチューニングを行い、ロバスト性を獲得することで分類精度を改善した。美的評価エンジンを Jetson TX1 に実装して構図推薦カメラを開発した。また、美的評価の高い写真が得られる方向へ構図を誘導し、撮影を指示する処理がリアルタイムに動作することを確認した。この構図推薦は初心者に見られるようなミス指摘したり、熟練者が意識するような構図を推薦して新たな気づきや構図を追求する機会を与えた。ユーザ評価によってデジタルカメラの金額的価値の上昇を確認した。

さらに撮影技術の習得効率を向上させるためには、美的評価に対する詳細な理由づけが必要と考えられる。今後の発展としては評価値の理由、および画像の部分領域ごとの改善点を明らかにすることを可能にしていきたい。

参考文献

[1] Murray, N., Marchesotti, L. and Perronnin, F.: AVA: A large-scale database for aesthetic visual analysis, *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2408–2415 (online), DOI: 10.1109/CVPR.2012.6247954 (2012).

[2] ブライアン・ピーターソン: ナショナルジオグラフィックプロの撮り方 構図を極める (ナショナル・ジオグラフィック), 日経ナショナルジオグラフィック社 (2013).

[3] 内池秀人, 福井麻衣子: 写真構図のルールブック, マイナビ (2012).

[4] 山田芳文: 写真は「構図」でよくなる! すぐに上達する厳選のテクニック 23, エムディエヌコーポレーション (2018).

[5] 家田 暁, 琴 智秀, 萩原将文: 感性を反映した構図修正による写真品質向上システム, 芸術科学会論文誌, Vol. 9, No. 4, pp. 163–172 (オンライン), DOI:

10.3756/artsci.9.163 (2010).

[6] 志津野之也, 濱川礼: 構図マッチング手法を用いた写真撮影時の自動構図決定手法, マルチメディア、分散協調とモバイルシンポジウム 2014 論文集, Vol. 2014, pp. 646–656 (2014).

[7] Bhattacharya, S., Sukthankar, R. and Shah, M.: A Framework for Photo-quality Assessment and Enhancement Based on Visual Aesthetics, *Proceedings of the 18th ACM International Conference on Multimedia*, MM '10, New York, NY, USA, ACM, pp. 271–280 (online), DOI: 10.1145/1873951.1873990 (2010).

[8] Lu, X., Lin, Z., Jin, H., Yang, J. and Wang, J. Z.: RAPID: Rating Pictorial Aesthetics Using Deep Learning, *Proceedings of the 22Nd ACM International Conference on Multimedia*, MM '14, New York, NY, USA, ACM, pp. 457–466 (online), DOI: 10.1145/2647868.2654927 (2014).

[9] Kao, Y., He, R. and Huang, K.: Deep Aesthetic Quality Assessment With Semantic Information, *IEEE Transactions on Image Processing*, Vol. 26, No. 3, pp. 1482–1495 (online), DOI: 10.1109/TIP.2017.2651399 (2017).

[10] Krizhevsky, A., Sutskever, I. and Hinton, G. E.: ImageNet Classification with Deep Convolutional Neural Networks, *Advances in Neural Information Processing Systems 25* (Pereira, F., Burges, C. J. C., Bottou, L. and Weinberger, K. Q., eds.), Curran Associates, Inc., pp. 1097–1105 (online), available from (<http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>) (2012).

[11] Yangqing J.: Caffe Model Zoo, http://caffe.berkeleyvision.org/model_zoo.html.

[12] Amazon Web Services, Inc.: 旧世代のインスタンス, <https://aws.amazon.com/jp/ec2/previous-generation/>.

[13] NVIDIA Corporation: NVIDIA JETSON, <https://www.nvidia.com/ja-jp/autonomous-machines/embedded-systems-dev-kits-modules/>.

[14] Connect Tech Inc.: Orbitty Carrier for NVIDIA Jetson TX2 & Jetson TX1, <http://connecttech.com/product/orbitty-carrier-for-nvidia-jetson-tx2-tx1/>.

[15] Sony Corporation, Sony Marketing Inc.: レンズスタイルカメラ ILCE-QX1, <https://www.sony.jp/ichigan/products/ILCE-QX1/>.

[16] SONY Corporation: Camera Remote API, <https://developer.sony.com/ja/develop/cameras/>.

付 録

A.1 構図の探索範囲

構図推薦は、一定のサイズの9枚の画像から最適な構図を探索している。より広い領域、狭い領域、あるいは回転させた領域、等の探索には現状では対応していないが、探索する画像を増やすことで可能になる。また、レンズの絞りについても現状は対応できていないが、絞りは画像の内容を大きく変える要素であり、撮影前に設定値を変更することで事前に画像を想定できるため、探索対象に含めることは可能である。しかしながら、いずれの場合も探索範囲を広げると計算量は増加する。