# 米飯画像の実寸推定に基づく面積を考慮したカロリー量推定

# 會下 拓実 柳井 啓司

†電気通信大学大学院情報理工学研究科 〒 182−8585 東京都調布市調布ヶ丘 1−5−1 E-mail: †ege-t@mm.inf.uec.ac.jp, ††yanai@cs.uec.ac.jp

あらまし 近年,健康志向の高まりにより様々な食事管理アプリケーションが開発されており,栄養学の知識のない人が外出先でも食事記録をすることが容易になっている.しかしこれらのアプリケーションでのカロリー量推定は,ユーザ入力による情報が必要であったり,栄養士を雇ったりと,人手のかかるものとなっている.食事画像からカロリー量などの栄養素の推定を自動で行うことができると健康管理の面で大きなメリットがあるが,現状では困難な問題である.一方,画像認識分野では Convolutional Neural Network(CNN) を用いた手法が主なタスクの最高精度を独占している.この CNN による食事画像の認識に関する研究も盛んに行われているが,高精度の食事画像からのカロリー量推定は実現されていない.そこで本研究では,面積を考慮した食事画像からのカロリー量推定を行う.そのために,まずCNN を用いた米飯画像からの実寸推定を行う.米飯粒は大きさが一定であるため,複数の米飯粒が密集した米飯画像から実すを直接推定する CNN を構築する.そして料理領域分割と実寸推定を組み合わせることで,面積を考慮した食事画像からのカロリー量推定を実現する.本手法では食事画像をテーブル面に垂直に真上から撮影することを仮定する.実験では,撮影した米飯画像に実寸情報をアノテーションすることで構築したデータセットを用いる.実寸推定の実験を行った結果,224 ピクセルあたりの実寸を推定したときの絶対誤差と相対誤差がそれぞれ 0.165cm と 6.394%となり,また,推定値と正解値の相関係数が 0.951 となり,高い相関が得られた.

キーワード カロリー量推定, 食事画像認識

### 1 はじめに

近年、健康志向の高まりにより様々な食事管理アプリケーションがリリースされ、栄養学の知識のない人がカロリー量を記録することが可能となっている。しかしそれらは複数の操作をユーザに要求し、リアルタイム性に欠けるものが多いため、より簡便な食事記録の方法が求められている。

そうした中で画像認識技術による食事画像からの自動認識が盛んに行われるようになっている。画像認識分野においてはDeep Convolutional Neural Network(CNN) により精度が飛躍的に向上しており、CNN による食事画像の認識に関する研究も盛んに行われ、実際にアプリケーションに利用された例も多い。

カロリー量は料理カテゴリおよび量に強く依存すると考えられ、料理カテゴリや量を食事画像から自動推定することが可能となれば、食事管理の面で有用である。食事画像からの料理カテゴリ分類においては既に CNN を用いた手法が高精度の分類を達成しており、最近では画像認識により食事画像から料理名の候補を自動で提案するアプリケーションも存在する。しかしカロリー量計算に不可欠な料理の量においては、料理カテゴリ毎に基準量を設けることで料理の量を考慮しないものも多い、料理量を考慮する場合であっても、ユーザ入力や基準物体が必要であるなど、料理の量の取得は非常に困難である。このように現状では、食事画像からのカロリー量推定は未解決の問題となっている。



図1 本研究において作成された米飯画像データセットの例. 各米飯画像に実寸情報が付与されている.

そこで本研究では面積を考慮した食事画像からのカロリー量推定を行う.そのために、まず CNN を用いた米飯画像からの実寸推定を行う.料理の量を考慮した食事画像からのカロリー量推定の既存研究では、対象の料理と同時に撮影された、大きさが既知の基準物体の領域と料理の領域を比較することで、料理の面積に基づいて料理のカロリー量を推定していた.本研究では、大きさが一定である米飯粒を大きさが既知の基準物体として扱い、複数の米飯粒が密集した米飯画像から実寸を直接推定する.米飯画像のパッチ画像を入力として、そのパッチ画像の一辺の長さの実寸を出力する CNN を構築する.そして料理検出と料理領域分割を行い、実寸推定と組み合わせることで面積を考慮した食事画像からのカロリー量推定を行う.ただし本手法では、食事画像をテーブル面に垂直に真上から撮影すること

を仮定する. 実験では撮影した米飯画像に実寸情報をアノテーションすることで構築したデータセット (図 1) を使用する.

まとめると、本研究では面積を考慮した食事画像からのカロリー量推定を行う。そのために、まず CNN を用いて米飯画像から実寸を推定する。米飯粒は大きさが一定であるため、米飯粒が密集した米飯画像から実寸を直接推定する CNN を構築する。そして料理領域分割と実寸推定を組み合わせることで、面積を考慮した食事画像からのカロリー量推定を行う。

# 2 関連研究

食事画像からのカロリー量推定にはいくつかのアプローチが存在するが、主要なアプローチは、推定された料理カテゴリと料理の面積や体積の情報から、事前に登録された料理カテゴリごとの単位面積当たりもしくは単位体積当たりのカロリー量の値を利用してカロリー量を推定する手法である.

Chen ら [1] は料理カテゴリを推定後, Kinect のような深度カメラにより料理の体積を推定し, 最終的にカロリー量を推定している. 深度カメラによる料理の体積の推定は正確であるが特殊なデバイスであるため, 一般の人が普段使用することは難しいと考えられる.

Kong ら [2] は Diet Cam という複数枚の画像からカロリー量を推定するアプリケーションを提案している. このアプリケーションは料理カテゴリ認識と領域分割を行い, さらに料理の三次元モデルの再構成を行い, 最終的に推定された体積の値からカロリー量を推定している. 三次元モデルの再構成では局所特徴量に基づくキーポイントマッチングとホモグラフィ推定が行われている. Dehais ら [3] の研究もこれに似ており, 皿の検出と領域分割, 料理カテゴリ分類を行い, 複数枚の画像から三次元モデルの再構成を行い, 最終的に炭水化物の量を推定している. このような複数視点からの画像により体積を推定する方法は, 事前にスマートフォンのカメラの較正を行わなくてはならなかったり, 正確に較正した地点から撮影を行わなくてはならず, ユーザに対する負担が大きいと考えられる.

Myers ら [4] は Im2Calories というアプリケーションを提案しており、食事/非食事の認識、複数品目の認識、深度推定、領域分割などの複数のタスクを CNN により行い、カロリー量を推定している。まず、食事/非食事認識により画像中に料理が存在するかを判定し、その後マルチラベル認識により画像中の複数の料理を認識する。次に深度推定と領域分割を行い、物体の三次元構造と料理の領域を抽出し、これらの情報を統合して料理の量を推定する。最後に料理カテゴリや量の情報から料理のカロリー量を推定している。この研究では、タスクごとに必要な学習データを独自に作成しているため、かなりのコストがかかると考えられる。また、カロリー量情報付きのデータセットが不足し、十分に評価が行われていない問題点がある。

Pouladzadeh ら [5] は料理とユーザの親指を同時に撮影することで指の大きさと比較を行い料理の大きさを求め,カロリー量を推定するシステムを提案している.しかし指の出し方や角

度, 映り方などによっては誤差が生じてしまう可能性がある. 本研究では, 大きさが一定である米飯粒を基準物体として実寸推定を行う.

岡元ら[6] は大きさが既知の基準物体と一緒に料理を撮影することで料理の体積を推定し、高精度のカロリー量推定を実現した。まず、基準物体と料理を一緒に撮影し、基準物体と料理のそれぞれの領域を抽出する。そして基準物体と料理の領域を比較して算出した料理の大きさからカロリー量を計算する。料理の領域の抽出では、まずエッジにより背景から皿領域を検出し、その皿領域に対して k-means により色情報に基づく領域分割を行い、最終的に GrabCut [7] により皿領域から料理領域を推定する。実験には基準物体と料理が一緒に写った画像が必要であるが、この手法は高精度の料理領域の推定を実現し、カロリー量推定では相対誤差 21%という精度を達成した。これに対して本研究では、米飯粒を基準物体として扱い、CNN を用いて米飯画像からの実寸推定を行う。

# 3 手 法

本研究では、料理の面積を考慮したカロリー量推定を行うために、料理検出と料理領域分割、そして米飯画像からの実寸推定を行い、食事画像から料理の実面積を推定する。提案手法では以下の流れで面積を考慮したカロリー量の推定を行う。

- (1) 対象の料理と米飯を同時に撮影
- (2) CNN を用いて各料理を検出
- (3) 各検出領域から CNN を用いて料理領域を抽出
- (4) 米飯画像から実寸を推定
- (5) 推定された実寸情報から各料理領域の実面積を推定
- (6) 推定された実面積に基づきカロリー量を推定

本手法では、CNN に基づく手法により、料理検出と料理領域分割を行い、さらに大きさが一定である米飯粒を基準物体として実寸を推定する CNN を構築する. 提案手法では、食事画像をテーブル面に垂直に料理の真上から撮影することを仮定する. 以下にそれぞれの処理の詳細を記述する.

#### 3.1 料理検出

本手法では各料理を検出するために Redmon らが提案した YOLOv2 [8] を使用する. YOLOv2 は以前に提案された CNN に基づく YOLO [9] を改善することで,高速かつ高精度な物体 検出を可能としている. 本手法では, YOLOv2 の学習には,大規模食事画像画像データセットである UECFood-100 [10] に含まれる画像 5500 枚に付与し直した料理バウンディングボックスを使用する. 画像 5000 枚を用いて YOLOv2 の学習を行い,500 枚での評価を行った結果, mAP (mean Average Precision)が 0.8 となり,高精度な料理検出を可能とした. 図 2 に,付与し直した料理バウンディングボックスを学習した YOLOv2 での料理検出の結果と, UECFood-100 を学習したクラス分類モデルにより各検出領域から推定された料理カテゴリを示す.

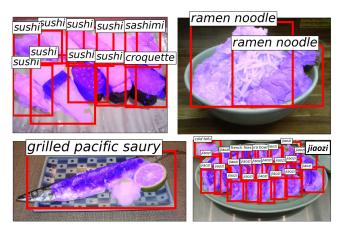


図 2 料理検出と料理領域分割の結果. 赤枠が推定されたバウンディン グボックス, 各赤枠上のタグが推定された料理カテゴリ, 赤い領 域が推定された料理領域である.

#### 3.2 料理領域分割

本手法では各料理の領域を推定するために、Ronneberger らが提案した U-Net [11] を使用する. U-Net は入力層付近の解像度の高い特徴マップを出力層付近に用いることで、高精度な領域分割を可能としている. 本手法では、U-Net の学習には、UECFood-100 [10] に含まれる画像 5500 枚に新たに付与した料理セグメンテーションマスクを使用する. 画像 5000 枚を用いて U-Net の学習を行い、500 枚での評価を行った結果、mIoU (mean Intersecting over Union) が 0.8 となり、高精度な料理領域分割を実現した. 図 2 に、新たに付与した料理セグメンテーションマスクを学習した U-Net での料理領域分割の結果を示す.

#### 3.3 米飯画像からの実寸推定

本手法では面積を考慮した食事画像からのカロリー量推定を行うために、CNNを用いて米飯画像からの実寸推定を行う.大きさが一定である米飯粒を基準物体として、米飯粒が密集する米飯画像から実寸を直接推定する CNN を構築する.図3に本手法の実寸推定モデルの構造を示す.

本研究で用いる CNN は VGG16 [12] に基づく. VGG16 は 畳み込み層が 13 層, 全結合層が 3 層の合計 16 層のネットワークであり, その性能と汎用性の高さから様々な研究に適用されている. 本手法の実寸推定モデルは図 3 のように, 実寸を出力する単一のユニットで構成される出力層を有する. 入力は米飯画像から得られるパッチ画像であり, 出力は入力されたパッチ画像の一辺の長さの実寸である. 本実験では入力パッチ画像のサイズを 224×224 とし, 出力は 224 ピクセルあたりの実寸となる.

次に米飯画像からの実寸推定での CNN の学習に用いる損失 関数について述べる. 実寸推定は回帰問題であるため, 本研究 では損失関数として式 (1) に示す 2 乗和誤差を使用する.

$$L = \frac{1}{n} \sum_{k=1}^{n} (x_k - y_k)^2 \tag{1}$$

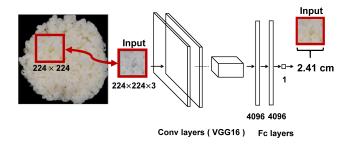


図 3 米飯画像から実寸を推定する CNN の構造.

### 3.4 カロリー量推定

本手法では、料理検出と料理領域分割により得られた料理領域情報と、米飯画像から推定された実寸情報から各料理の実面積を求め、そこから岡元ら[6]の手法に従い、各料理のカロリー量を推定する。岡元らは事前に複数サイズの料理のカロリー量より学習した回帰曲線に基づき、実面積情報からカロリー量を推定した。料理カテゴリ毎に学習された回帰曲線を用いることで、牛丼やごはんといった深さのある皿に盛られることが多い料理に関しても2次曲線のようなフィッティングを行うことが可能である。岡元らは、学習用の画像として1食品あたり3サイズの基準物体と共に写ったカロリー量情報付き画像を用意し、その画像から求められた面積と付与されたカロリー量情報から回帰曲線を求めた。また、対象食品として日本食を中心とした20種類が選択された。本手法では、岡元ら[6]が学習した回帰曲線に基づき、各料理の実面積情報からカロリー量の推定を行う。

# 4 実寸情報付き米飯画像データセット

本研究では米飯画像からの実寸推定を行うために、実寸情報付き米飯画像データセットを新たに作成した。まず米飯画像の撮影では、図4のように2種類のカメラ (COOLPIX AW120とiPhone8 Plus)を使用し、各水量(米150gに対して水180ml、200ml、220ml)を用いて炊飯したそれぞれの米飯の撮影を行った。カメラの種類と炊飯時の水量の組み合わせごとに60枚撮影したため、全体で360枚の画像が収集された。また、1枚撮影する毎に被写体の米飯とカメラとの間の距離を変更し、5枚撮影する毎に米飯の盛り付けを変更することで、多様な米飯画像を収集した。

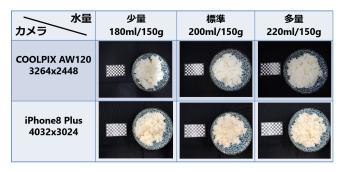


図 4 本研究において撮影された米飯画像の例.

次に撮影した米飯画像への実寸情報の付与を行った. 米飯画

像の撮影時に使用した茶碗の直径に基づき、1 ピクセルあたりの実寸を算出した. さらに背景情報を除去するためにセグメンテーションマスクの付与を行った. 図 1 に本研究において作成された米飯画像データセットを示す.

## 5 実 験

本実験では、米飯画像からの実寸推定と面積を考慮した食事画像からのカロリー量推定を行う.米飯画像からの実寸推定での学習と評価には、本研究で作成された図1の実寸情報付き米飯画像データセットを用いる.面積を考慮したカロリー量推定では、料理領域分割と実寸推定を組み合わせることで、実面積を推定し、そこからカロリー量を推定する.テスト画像として、UECFood-100[10]の複数品料理画像のうち、画像に米飯が含まれるものを用いる.

#### 5.1 米飯画像からの実寸推定

本実験では米飯画像からの実寸推定を行う. データセットとして本研究で収集した実寸情報付き米飯画像データセットを用いる. カメラの種類と炊飯時の水量の組み合わせに基づきデータセットを6つに分割し,5つを学習に使用し,残りの1つを評価に使用する. そのため学習画像と評価画像はそれぞれ300枚と60枚となる.

学習時には、1 枚の米飯画像のランダムな位置から切り抜かれた 1 枚の  $224 \times 224$  のパッチ画像に対して推定された実寸に対して学習が行われる。また、前処理として米飯画像の拡大縮小と左右反転、回転を行う。拡大縮小では、ある米飯画像をn 倍に

拡大縮小した場合,その米飯画像に付与された実寸情報を $\frac{1}{n}$ 倍にする.評価時には,1枚の米飯画像から  $4\times 4$  のグリッドサンプリングによって切り抜かれた 16 枚のパッチ画像それぞれから推定された実寸の平均値が最終的な出力となる.さらに学習時と評価時の両方において,背景領域が 50%以上を占めるパッチ画像は背景画像として除去される.図 5 と図 6 に学習時と評価時それぞれの入力パッチ画像を示す.



図 5 学習時における入力パッチ画像群 (バッチサイズ:16). 16 枚の各 米飯画像から 1 枚のパッチ画像がランダムな位置から切り抜か れる. 背景画像は除去される.

本実験において米飯画像からの実寸推定に用いる CNN の構造は VGG16 [12] に基づく. 出力層以外の層では, ImageNet の 1000 種類分類タスクにより学習済みの重みを学習時の初期値として使用する. 最適化手法として SGD を使用し, momentum値を 0.9 する. バッチサイズを 16 として学習率  $10^{-5}$  において約 1,900 回反復する.



図 6 評価時における入力パッチ画像群. 1 枚の米飯画像から 4×4 の グリッドサンプリングによって切り抜かれる. 背景画像は除去されるため 14 枚となっている.

評価指標として絶対誤差,相対誤差,推定値と正解値の相関係数,相対誤差 5%,10%,20%以内の推定値の割合を用いる. 絶対誤差は推定値と正解値の差の絶対値であり,相対誤差は正解値に対する絶対誤差の割合である. なお,評価には224 ピクセルあたりの実寸における推定値と正解値を用いる.

表 1 に実寸推定の結果を示し、図 7 に推定値と正解値の相関を示す。224 ピクセルあたりの実寸を推定したときの平均絶対誤差と平均相対誤差がそれぞれ 0.165cm と 6.394%となり、また、推定値と正解値の平均相関係数が 0.951 となり、高い相関が得られた。 すべての学習データと評価データの組み合わせにお

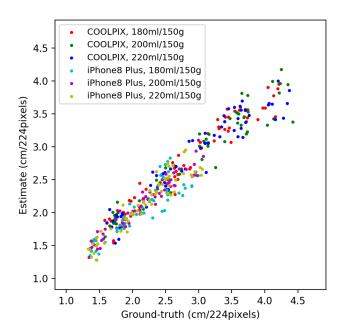


図7 米飯画像からの実寸推定における推定値と正解値の相関関係.

いて、相対誤差が 10%より小さく、相関係数が 0.9 より大きくなり、また、ほとんどの推定値が相対誤差 20%以内 (%) に含まれている. このような結果から、大きさが既知の基準物体として米飯粒を用いることが有効であることが示された. また本研究では、米飯粒が密集した状態の米飯画像から実寸を直接推定する CNN を学習することで、米飯粒ひと粒ひと粒の大きさや向きのばらつきに頑健な推定を可能にしたと考えられる.

# 5.2 面積を考慮したカロリー量推定

本実験では CNN による料理領域分割と米飯画像からの実寸推定を組み合わせることで, 領域分割に基づく面積を考慮した食事画像からのカロリー量推定を行う. YOLOv2 [8] での料理検

表 1 米飯画像からの実寸推定 (224 ピクセルあたりの実寸における評価)

評価データ	AP 計画 元 ( )	和补细 关 (04)	+口目目 125 米h	#마찬[레 꽃 ㅌ여 이 라 (여)	10 사례 중 10 에 다 다 (여)	和計劃表 2000 17中 (04)
評価テータ	絶対誤差 (cm)	相対誤差 (%)	相関係数	相対誤差 5%以内 (%)	相対誤差 10%以内 (%)	相対誤差 20%以内 (%)
COOLPIX,180ml/150g	0.212	7.182	0.958	41.667	75.000	91.667
$\rm COOLPIX,200ml/150g$	0.178	6.550	0.973	43.333	76.667	93.333
$\rm COOLPIX,220ml/150g$	0.197	6.668	0.962	48.333	78.333	90.000
iPhone8 Plus, 180ml/150g	0.127	5.652	0.945	50.000	75.000	98.333
i Phone 8 Plus, $200\mathrm{ml}/150\mathrm{g}$	0.170	7.512	0.903	43.333	68.333	88.333
i P none 8 Plus, $220\mathrm{ml}/150\mathrm{g}$	0.105	4.800	0.967	58.333	88.333	98.333

出により得られた各料理領域に対して、クラス分類と U-Net [11] による領域分割を行うことで各料理のカテゴリと領域を推定する. さらに抽出した米飯画像から実寸を推定することで得られた 1 ピクセル当たりの実面積を基準として、料理の実面積を算出し、そこから岡元ら [6] の手法に従い、カロリー量を推定する. テスト画像として、UECFood-100 [10] の複数品料理画像のうち、画像に米飯が含まれるものを用いる. 図 8 に実面積推定とカロリー量推定の結果を示す.

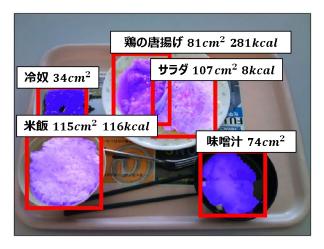


図 8 米飯を含む複数品料理画像からの面積を考慮したカロリー量推定の結果.各赤枠上のタグが推定された実面積とカロリー量(対象食品20種類のみ)である.実際にはテーブル面に垂直に真上から撮影することを想定している.

## 6 おわりに

本研究では料理面積を考慮した食事画像からのカロリー量推定を行うために、CNNを用いて米飯画像からの実寸推定を行った。米飯粒は大きさが一定であるため、複数の米飯粒が密集した米飯画像から実寸を直接推定する CNNを構築した。実験では224ピクセルあたりの実寸を推定したときの平均絶対誤差と平均相対誤差がそれぞれ0.165cmと6.394%となり、また、推定値と正解値の平均相関係数が0.951となり、高い相関が得られた。そして料理検出と料理領域分割を行い、米飯画像からの実寸推定と組み合わせることで、面積を考慮した食事画像からのカロリー量推定を行った。

今後の課題として,面積を考慮したカロリー量推定の評価を 行うことがある.評価を行うために,米飯が含まれるカロリー 量情報付き複数品料理画像データセットを作成する予定である. また、米飯がない状況においても実面積推定が可能なシステムを構築するために、領域分割と基準物体を用いる手法[6],[13] や、iPhone の深度対応カメラから得られる深度情報と組み合わせることも検討している.

#### 文 献

- [1] M. Chen, Y. Yang, C. Ho, S. Wang, E. Liu, E. Chang, C. Yeh, and M. Ouhyoung. Automatic chinese food identification and quantity estimation. In Proc. of SIGGRAPH Asia Technical Briefs, pp. 1–4, 2012.
- [2] F. Kong and J. Tan. Dietcam: Automatic dietary assessment with mobile camera phones. *Pervasive Mob. Comput.*, Vol. 8, No. 1, pp. 147–163, 2012.
- [3] J. Dehais, M. Anthimopoulos, and S. Mougiakakou. Gocarb: A smartphone application for automatic assessment of carbohydrate intake. In Proc. of ACM MM Workshop on Multimedia Assisted Dietary Management, 2016.
- [4] A. Myers, N. Johnston, V. Rathod, A. Korattikara, A. Gorban, N. Silberman, S. Guadarrama, G. Papandreou, J. Huang, and P. K. Murphy. Im2calories: towards an automated mobile vision food diary. In *Proc. of IEEE Inter*national Conference on Computer Vision, pp. 1233–1241, 2015
- [5] P. Pouladzadeh, S. Shirmohammadi, and R. Almaghrabi. Measuring calorie and nutrition from food image. In *IEEE Transactions on Instrumentation and Measurement*, pp. 1947–1956, 2014.
- [6] K. Okamoto and K. Yanai. An automatic calorie estimation system of food images on a smartphone. In Proc. of ACM MM Workshop on Multimedia Assisted Dietary Management, 2016.
- [7] C. Rother, V. Kolmogorov, and A. Blake. "GrabCut": Interactive foreground extraction using iterated graph cuts. ACM Trans. Graph., Vol. 23, No. 3, pp. 309–314, 2004.
- [8] J. Redmon and A. Farhadi. YOLO9000: Better, faster, stronger. In Proc. of IEEE Computer Vision and Pattern Recognition, 2017.
- [9] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You Only Look Once: Unified, real-time object detection. In Proc. of IEEE Computer Vision and Pattern Recognition, 2016
- [10] Y. Matsuda, H. Hajime, and K. Yanai. Recognition of multiple-food images by detecting candidate regions. In Proc. of IEEE International Conference on Multimedia and Expo, pp. 25–30, 2012.
- [11] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. Springer, pp. 234–241, 2015.
- [12] K. Simonyan and A Zisserman. Very deep convolutional networks for large-scale image recognition. In arXiv preprint arXiv:1409.1556, 2014.
- [13] W. Shimoda and K. Yanai. CNN-based food image segmentation without pixel-wise annotation. In Proc. of IAPR International Conference on Image Analysis and Processing, 2015.