

コンテキスト情報に基づく スマートグラスを用いた家電選択手法の提案

孔全^{1,2} 前川卓也^{1,2,a)} 宮西大樹² 須山敬之²

受付日 2016年12月10日, 採録日 2017年7月4日

概要: 本論文では, Google Glass に代表されるようなスマートグラスを用いて, ホームネットワークに接続したスマートホーム内の家電を容易に選択する手法の提案およびそれを実現するための手法設計を行う. 提案手法では, Google Glass に搭載されたユーザの視線方向を撮影するカメラを用いて, ユーザが注目している家電を特定することで, 家電の選択を実現する. すなわち, ユーザは頭部を注目する家電の方向に向けるだけで, その家電を選択できる. 本研究では, ユーザの家電への注目行為をスマートグラス上の方位センサデータにより検出する. さらに操作したい家電の特定を実現するため, 本研究の特徴として, カメラに加えて加速度センサなどの他のセンサから得られた行動や位置に関するコンテキスト情報を用いて, 高精度な家電の特定を実現する点があげられる. これは, ユーザが操作したいと考えている家電は, そのユーザのコンテキスト情報に強く関係していると考えられるためである. 提案手法では, ディープラーニング技術を用いて画像特徴を抽出し, ノンパラメトリッククラスタリングを用いたラベルなしセンサデータからの行動や位置に関連する特徴抽出を行ったあと, マルチカーネル学習を用いてそれらの特徴を組み合わせた家電選択を行う. 評価実験では, 実際に様々なネットワーク家電が設置されたスマートホームでセンサデータを取得し, 提案手法の効果を検証した.

キーワード: ウェアラブルコンピューティング, スマートグラス, 電化製品, コンテキスト情報

Contextual Information Based Home Appliances Selection with Smart Glass

QUAN KONG^{1,2} TAKUYA MAEKAWA^{1,2,a)} TAIKI MIYANISHI² TAKAYUKI SUYAMA²

Received: December 10, 2016, Accepted: July 4, 2017

Abstract: We propose a method for selecting home appliances using a smart glass, which facilitates the control of network-connected appliances in a smart house. In this paper, we mainly describe on how we achieve a reliable performance on appliance selection by using a smart glass. Our proposed method is image-based appliance selection and enables smart glass users to easily select a particular appliance by just looking at it. The main feature of our method is that it achieves high precision appliance selection using user contextual information such as position and activity, inferred from various sensor data in addition to camera images captured by the glass because such contextual information is greatly related in the home appliance that a user wants to control in her daily life. Our experimental results, which use sensor data obtained in an actual house equipped with many network-connected appliances, show the effectiveness of our method.

Keywords: wearable computers, smart glass, home appliances, context information

¹ 大阪大学大学院情報科学研究科
Graduate School of Information Science and Technology,
Osaka University, Suita, Osaka 565-0871, Japan

² 国際電気通信基礎技術研究所
Advanced Telecommunications Research Institute International (ATR), Souraku-gun, Kyoto 619-0237, Japan

a) maekawa@ist.osaka-u.ac.jp

1. はじめに

近年, 様々な電化製品をホームネットワークに接続することにより, それらの稼働状況やエネルギー消費量などを簡単に取得することが可能になりつつある. また, 家電

情報を取得することに加えて、ホームネットワークを通じて、家電をコントロールすることも可能になりつつある。たとえば、ネットワークを通じて、ウェアラブルデバイスやタブレット端末により家庭内の照明システムを操作することがその一例である。特にウェアラブルデバイスにより家電を制御することには次のようなメリットがある：(1) リモコンが手の届くところになくとも家電を操作できる。(2) ハンズフリーの操作を実現すれば（たとえば、音声操作など）、他の行動によって両手が塞がっている場合でも家電の操作ができる。また、身障者に対しても、ハンズフリーの操作方法は有効なアプローチと考えられる。

ユーザが家電を操作したいとき、まず家にある数多くの家電から操作したい家電を1つ選択する必要がある。ウェアラブルやユビキタスコンピューティングの研究分野では、家電製品を選択するためのいくつかの一般的な方法がこれまでに研究されている [4], [22], [23]。たとえば、音声、カメラ画像、ジェスチャおよびビーコンなどに基づく方法がこれまで研究されてきた。しかし、たとえば、ビーコンを用いた方法はビーコンを家電に取り付ける必要がある。また、音声を用いた方法では、長い言葉で操作したい家電を特定する必要がある。そして、ジェスチャを用いた方法ではユーザが家電ごとに対応するジェスチャを覚える必要がある。

本研究では、ウェアラブルデバイスの1つであるスマートグラスを使用して、家電を選択する方法の提案およびそれを実現するための手法設計を行う。提案手法では、画像に基づく家電の選択手法として、スマートグラスの1つである Google Glass を装着するユーザが頭を選択したい家電に向けるだけで、その際のセンサデータを用いて家電を特定する。それを実現するため、提案手法では、まずユーザの家電への注目行為を検出し、そして、注目時刻における画像からユーザが選択したい家電を認識する。Google Glass は頭の向きの方向を撮影するカメラを備えており、そのカメラを用いて選択したい家電をとらえる。しかし、画像のみを用いるアプローチには、似た外見を持つ家電に対して、それらを区別する画像特徴が少ないため、認識が困難となるような問題が存在する。本研究では、画像に基づく高精度な家電選択手法を実現するため、カメラの画像に加えて加速度センサなど他のセンサから得られたコンテキスト情報を利用する。本研究でのコンテキスト情報とは Google Glass を装着するユーザの屋内位置情報や行動情報など、ユーザが操作したいと考えている家電に強く関係しているコンテキスト情報のことである。たとえば、キッチンで料理をしている際は、キッチン家電を操作したいと考えることが多いと考える。また、コンテキスト情報は、画像のみでは区別できない家電に対しても効果的である。たとえば、寝室にあるエアコンとリビングにある外見が同じエアコンをユーザの位置情報を付加的に用いて区別でき

る。本研究では、屋内位置や行動の種類を教師なしの学習方法により認識するため、ユーザの位置と行動を推定するためのラベル（位置の場合は寝室、リビングなど、行動の場合は調理、寝るなど）付きトレーニングデータを用意する必要はない。

本論文の以降の構成は以下のとおりである。まず家電選択・操作の関連研究を紹介したあと、Google Glass による家電選択システムについて述べる。そして、カメラ画像に加えてコンテキスト情報を用いる家電選択手法の詳細を説明する。また、実際の環境で収集したセンサデータを用いて提案手法を評価する。最後に本論文のまとめを行う。

2. 関連研究

家電の選択および操作手法は大きく分けると、主に音声、ジェスチャ、ビーコン、視線およびビジョンに基づくアプローチがある。

2.1 音声

音声に基づく方法がこれまでに開発されている。たとえば、Christensen ら [4] は音声で家電を操作するクラウドベースのシステムを開発した。しかし、曖昧性を排除して家電を正確に特定するためには長い家電の名前を読み上げる必要がある。また、日常生活の騒音からの影響も受けやすい。

2.2 ジェスチャ

ジェスチャに基づくアプローチは、選択したい家電を特定するため、身体装着型センサを利用している [22]。たとえば、ハンドジェスチャを認識するために、身体装着型カメラや手に装着する慣性センサが多く使われている。このアプローチでは、各家電ごとに1つのジェスチャが関連付けられているため、家の中にある多くの電化製品に対して、その関連付けを思い出すのは難しい。

2.3 ビーコン

ビーコンに基づくアプローチは電化製品に赤外線受信機やLED ビーコンを取り付ける必要があり、導入コストとメンテナンスコストが大きくなる [19], [23]。一方、本提案手法は電化製品に外部装置を取り付ける必要はない。

2.4 画像と視線

画像に基づくアプローチは主にウェアラブルやスマートフォンカメラを用いて操作対象を特定する [17], [18], [20]。しかし、家電に添付する AR マーカ必要とする手法もある。また、視線に基づくアプローチでは、ユーザの視線方向を撮影するカメラの画像を用いて、家電を特定する [5], [21]。このアプローチを用いる研究では実際に目の動きを検知することによって視線方向を検出しているものもある。そ

して、検出された視線と画像を用いて注目している家電を特定する。しかし、画像や視線ベースの手法は視覚的特徴が少ない、外見が似た家電を区別するのは困難であり、さらに暗い環境では動作しない可能性もある。本提案手法はユーザの視線方向（頭の向いている）を撮影した画像を用いたアプローチであるが、その画像に加えてユーザのコンテキスト情報を用いて家電選択の性能を向上させることを特徴とする。

2.5 既存手法におけるコンテキスト情報の適用可能性

本提案のコンテキスト情報を用いる方法は他の家電選択手法にも適用可能と考える。たとえば、音声に基づく家電選択手法に対して、ユーザの位置情報を事前知識として与えることで、「テレビ」などのあいまいな入力音声で家電を選択することも可能である。

3. 基本システム

3.1 想定環境

本研究では、ユーザは Google Glass など多種のセンサを搭載したスマートグラスを装着すること、また、ユーザは Google Glass とペアリングされているスマートフォンを持っていることを想定している。Google Glass は家のホームネットワークに接続され、また、テレビ、エアコン、ライトなどの電化製品も同じネットワークに接続されている。ユーザは、ネットワークに接続した Google Glass を利用して、家電の選択および選択した家電の操作を行う。スマートフォンはコンテキスト情報に関わるセンサデータ（加速度、Wi-Fi 信号、音声など）を計測する。

3.2 アプリケーション

次に、家電の選択および操作を行うための Google Glass（以降では Glass と呼ぶ）にインストールするアプリケーションについて説明する。図 1 は Glass のディスプレイに表示された情報をユーザが見ている実世界に重畳したイメージを示す。Glass のユーザインタフェースではユーザに提示されるひとまとまりの情報はカードと呼ばれる単位で管理される。本研究では各家電ごとに 1 枚のカードを用意している。このカードは Glass のディスプレイに表示さ

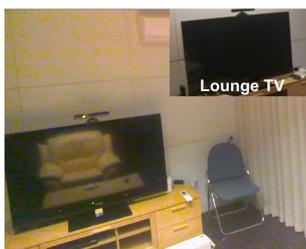


図 1 アプリケーションを利用するユーザの視野のイメージ
Fig. 1 Example view of our application user.

れる。ユーザは Glass の側面にあるタッチパッドをスワイプして、操作したい家電に関するカードを探す。家電、すなわちカードを選択した後、ユーザは選択した家電に関する操作を Glass のタッチパッドにより行うことができる。たとえば、エアコンの ON/OFF などの操作を行う。本研究はこのようなアプリケーションを基に、ユーザが操作したい家電に頭を向けるだけで、家電を選択する手法を実現する。すなわち、タッチパッドで操作したい家電のカードを探す代わりに、カメラ画像とコンテキスト情報を表すセンサデータを用いて推定された家電に対応するカードを Glass のディスプレイに自動的に表示する。そして、表示されたカードが正しい場合、すなわち推定された家電が実際にユーザが操作したい家電である場合、ユーザは Glass にあるタッチパッドで家電の操作を継続して行う。また、推定結果が正しくない場合、ユーザは手動でタッチパッドにより操作したい家電に対応するカードを選択する。さらに、推定結果が正しくない場合に、再びカードを手動で選択する負担を低減するため、Glass に表示するカードの並びを推定結果の尤もらしさのランキングによりソートする。このような方式を用いることで、ラベル付きの画像データも低負担で取得することが可能になる。

3.3 システム導入

システムを家庭に導入する際、ホームネットワークに接続された家電の MAC アドレスは事前に Glass の家電選択アプリケーションに登録されていると想定する。アプリケーションの 1 枚のカードは 1 台の家電の画像とその家電の名前（たとえば Lounge TV）からなる（図 1）。ユーザは最初に家電の前で家電の画像を撮り、その名前を決め、それを MAC アドレスと関連付けるカードの登録作業を行う必要がある。登録作業は以下の手順で行う。ステップ 1：ユーザは登録したい家電の前で、アプリケーションにあらかじめ用意された家電の MAC アドレスや出荷時のデフォルトの家電名を表示するカード群から登録したい家電に対応するであろうカードを選択する。選択されたカードから家電をコントロール可能（ON/OFF など）な場合、次のステップに進む。ステップ 2：ユーザは家電の前で複数枚家電の画像（実験では 30FPS の動画を約 10 秒）を撮り、その後、サーバに転送し、サーバ上でビデオの各フレームが静止画として抽出され、初期の学習データとして使用される。ステップ 3：ユーザは、家電の名前をデフォルトの名前から変更し、アプリケーションはその名前とビデオのサムネイル画像 1 枚をそのカードの新しい内容として差し替える。登録作業後、ユーザが日常生活でシステムを繰り返し使用することでラベル付きの画像とセンサデータが蓄積されていく。このトレーニングデータを用いてモデルのパラメータを逐次更新していく。また、アプリケーションはバックグラウンドで実行されるため、Glass のカメラは常

時オンにする必要はなく、ユーザの家電への注視が検出されたときのみカメラが自動的に起動される仕組みとなる。

4. 提案手法

4.1 概要

本研究では、Glass 上のカメラ、照度センサ、方位センサおよびユーザが携帯するスマートフォン上の加速度センサ、マイク、Wi-Fi モジュールをセンサとして利用する。カメラは視線（頭）方向の画像をとらえるために、方位センサは視線方向の計測に利用する。また、加速度センサ、マイクおよび照度センサは利用者の行動情報を推定するために利用し、Wi-Fi モジュールは環境内に設置した複数の Wi-Fi アクセスポイント（AP）から受信した信号強度情報を用いてユーザの屋内位置を推定するために用いる。

図 2 に提案手法の概要図を示す。提案手法では、まず Glass に搭載された方位センサを用いて、ユーザが何らかの対象に対して注視を行ったかどうかを検出する。注視の動作を行っていた場合、その他のセンサデータも用いて、実際にどの家電を注視したかを特定する。利用者の行動および位置情報の推定は教師なしクラスタリング手法の 1 つである Infinite Gaussian mixture model (IGMM) を用いることにより行う。最後に、推定された行動や位置情報の確率と Deep convolutional neural network (DCNN) から抽出された画像特徴を、Multiple kernel learning (MKL) を基にした家電選択モデルの入力として用いる。

4.2 注視の検出

家電がユーザの顔の正面にない場合、ユーザは顔を家電の方向に向け、家電に注目する。図 3 は実験参加者がキッチンでの調理作業終了後、リビングルームに移動し、エアコンとテレビをオンにしたときに、Glass から得られた時系列の方位センサデータの例を示している。また、図 3 はこれらの行動中に Glass により撮影された画像も示している。赤、緑、青の線は、それぞれ、 x 軸、 y 軸と z 軸のデータを示す。家電に視線を移すために頭部を回転させること

により、方位センサデータが大きく変化しあつたあと、その家電を注視するため頭部を固定することにより、方位センサデータには変化が見られなくなる。そこで、そのセンサデータが大きく変化し直後（1 秒以内）に存在する、変化が少ない静的なセンサデータセグメントを発見することにより、注視の検出を行う。本研究では、スライディングウィンドウを用いて、センサデータの分散の変化が大きいセグメントと静的なセンサデータセグメントを発見する。ウィンドウ内の 3 軸センサデータの分散は下の式により計算する：

$$v = \frac{1}{T} \sum_{i=t}^{t+T-1} (\bar{x} - x_i)^2 + \frac{1}{T} \sum_{i=t}^{t+T-1} (\bar{y} - y_i)^2 + \frac{1}{T} \sum_{i=t}^{t+T-1} (\bar{z} - z_i)^2$$

ここで、 T はウィンドウサイズ、 t はウィンドウ内の最初のデータサンプルのインデックスとなる。

計算された v が閾値より低い場合、ウィンドウ（データセグメント）は静的なものに見なし、そうでない場合、変化が大きいデータセグメントとする。

4.3 Deep Convolutional Neural Networks による画像特徴の抽出

注視を検出した後、その時刻においてカメラがとらえた画像データから特徴を抽出する。抽出された特徴は家電選択モデルの入力となる。注視が時刻 t_s に発生したとすると、 t_s から $t_s + 1$ [秒] の間に撮影された画像から特徴を抽出する。特徴の抽出は、一般的な物体の認識に強い DCNN により行う。本研究で用いた DCNN の入力画像は $1,280 \times 720$ ピクセルである。利用した DCNN は ILSVRC-2012 の画像データセットで事前訓練済みの 7 層のニューラルネットとなる [14]。6 層までが畳み込み層（convolution layer）で、その layer-6 の活性化信号（出力結果）を特徴量（4,096 次

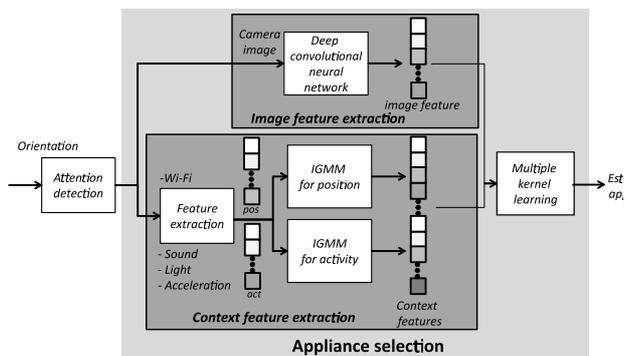


図 2 提案手法の概要

Fig. 2 Overview of our proposed method.

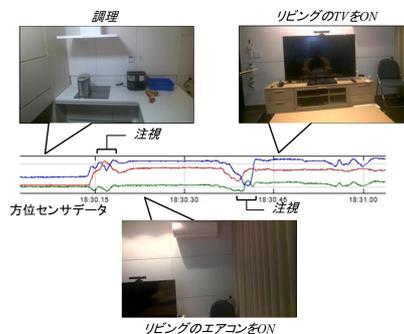


図 3 Google Glass から計測した時系列の方位センサデータの例。赤、緑、青の線は、それぞれ、 x 軸、 y 軸と z 軸のデータを示す

Fig. 3 Example of time series orientation data obtained from smart glass. Red, green and blue line shows x , y , and z axis respectively.

元のベクトル)として扱う。layer-7は語彙表現層となるため、layer-6の出力を利用する[8]。

DCNNによる画像特徴の抽出には高い計算能力が必要とされるため、本研究では、画像特徴の抽出はホームネットワークに接続した高性能サーバで行う。画像以外のセンサデータの特徴の抽出はスマートフォンで行う。

4.4 IGMMの入力となる特徴の抽出

ここで、各センサが計測したデータから特徴を抽出する方法を述べる。これらの特徴はIGMMの入力となり、抽出された特徴はユーザの行動および屋内位置の推定に利用される。注視(静的なウィンドウ)が検出された時刻を t_s とし、それに対応する変化が大きいウィンドウが検出された時刻を t_l ($t_l < t_s$)とする。

・**加速度センサ(スマートフォン)**:これまでの多くの行動認識の研究では、たとえば歩く、立つ、寝るなどの行動の認識に加速度データを利用している[16]。3軸方向(x, y, z)の加速度信号の変化はスマートフォン装着位置(ズボンのポケットや胸のポケット)や向きに大きく影響されてしまうため、従来の研究[10]で提案された加速度合成信号 $R_i = \arcsin\left(\frac{z_i}{\sqrt{x_i^2 + y_i^2 + z_i^2}}\right)$ を利用する。ここで、 R_i は i 番目の加速度合成信号となる。本研究では、時刻 $t_l - w$ から時刻 t_l までの合成信号データから平均と分散を計算し、加速度の特徴量として利用する。

・**マイク(スマートフォン)**:音声データも多くの行動認識研究で利用されている[15]。音声データは日常生活の中で音に関わる行動、たとえば、水を流す、会話や入浴などの認識に用いられる。メル周波数ケプストラム(MFCC)は環境音認識の研究において最も性能の高い特徴抽出方法といわれており[6]、本研究では、時刻 $t_l - w$ から時刻 t_l までの音声データから抽出したMFCCコンポーネントの各係数の平均を特徴量として用いる。

・**照度センサ(Google Glass)**:照度データはユーザが明るいところにいるかどうかの検出に有用である。本研究では、時刻 t_l から時刻 t_s までの照度データから計算した平均値を照度の特徴量とする。

・**Wi-Fiモジュール(スマートフォン)**:多くの屋内位置推定システム研究では、Wi-Fi信号強度を用いている[12]。本研究では、時刻 t_s に計測されたそれぞれのアクセスポイントのWi-Fi信号強度をWi-Fiの特徴量とする。

・**方位センサ(Google Glass)**:方位データはユーザの顔(視線)の向きを示すデータとなり、ユーザがどの家電に注目しているかに強く関係している。本研究では、時刻 t_s に計測された方位データの各軸(x, y, z)の平均値を特徴量とする。

4.5 教師なし行動認識および位置の推定

行動と屋内位置の推定には教師なし学習のアプローチを用いる。推定された行動と位置情報(の確率)は家電選択モデルの入力として用いられる。加速度、照度、マイクの特徴量は行動認識モデルの入力とする。そして、Wi-Fiの特徴量は位置推定モデルの入力として用いられる。図2に示したように、時刻 t の行動認識モデルに関連する加速度、照度、マイクの特徴量を連結し、 $S_{act,t}$ とする。また、時刻 t の位置推定モデルに関連するWi-Fiの特徴量を連結し、 $S_{pos,t}$ とする。連結した特徴量ベクトルを用いて以下の方法で行動認識と位置推定モデルを学習する。

・**学習フェーズ**:ユーザが日常生活で収集したラベルなしのセンサデータから抽出した特徴ベクトルを用いて、行動および位置モデルを学習する。本研究では、混合ガウス分布を用いて行動および屋内位置をモデル化する。1つのガウス分布は1つの行動や位置パターンのクラスタに対応し、たとえば、行動の場合は「調理」や「睡眠」などのクラスタに、位置の場合は「キッチン」や「寝室」などのクラスタに対応する。しかし行動と位置パターンの数は事前に既知ではなく、ユーザによっても異なるため、本研究はノンパラメトリックベイズ手法を用いて自動的に、クラスタ数を決定する。教師なし学習のため、行動や位置モデル内のクラスタには「調理」や「寝室」などようなラベルではなく、単に“activity cluster 01”などのラベルが付与される。

具体的には、IGMMを用いて、ノンパラメトリック教師なしクラスタリングを行う。IGMMはディリクレ過程のガウス混合モデル(GMM)であり[9]、GMMの混合数を無限大とおくことにより、混合数を事前に与える問題を回避している[3]。従来の混合モデルと異なり、IGMMはモデルの推論過程の中でクラスタ数も同時に推定する。すなわちIGMMではデータによりよく適応するよう、クラスタ数を必要に応じて任意に増加させることができる。詳細は文献[3]を参照されたい。

・**テストフェーズ**:新しいテスト用特徴ベクトルが得られたとき、テストベクトルと訓練したIGMMにおける各クラスタ間の距離を計算することで、そのベクトルがどのクラスタに属するか、すなわち、どの行動を行っているかもしくはどの位置にいるかを求める。具体的には、各クラスタとの距離の逆数からなるベクトルをコンテキスト特徴量とし、IGMMからの出力とする。そして、図2に示すように、行動認識用および位置推定用IGMMの出力を家電選択モデルの入力とする。

4.6 家電の選択

4.6.1 概要

操作したい家電はカメラの画像とコンテキスト情報により推定される。すなわち家電選択モデルの入力は、DCNNにより抽出した画像特徴と、IGMMが出力した距離(コ

ンテキスト特徴量)である。本研究では、異なるソースの特徴から構成されるデータにうまく対処することができる Multiple kernel learning (MKL) を基にした家電選択モデルを実現する。MKL は異なる基底カーネルの線形結合により構成され、本研究では異なるデータソースごと (画像とコンテキスト) に特化したカーネルを用意する。MKL を用いて入力ベクトルを分類し、家電のクラスを出力する。

4.6.2 マルチカーネルによる家電選択

カーネル関数はインスタンス間の距離を計算する関数であり、特徴空間における線形識別関数を定めるために用いられる。N 個のトレーニングインスタンス (ベクトル) $\{\mathbf{x}_i \in \mathcal{X}\}_{i=1}^N$ から、未知のテストインスタンス \mathbf{x}_* の推定に用いられる識別関数は以下のように定義される。

$$f(\mathbf{x}_*) = \mathbf{a}^T \mathbf{k}_* + b$$

ここで、a と b は各トレーニングインスタンスに対するベクトルの重みとバイアスとなる。また、 $\mathbf{k}_* = [k(x_1, \mathbf{x}_*) \dots k(x_N, \mathbf{x}_*)]^T$ であり、 $k(\cdot, \cdot)$ は2つのインスタンス間の距離 (非類似度) を計算するカーネル関数である。

MKL を用いた家電選択モデルにおける家電分類の識別関数として、画像とコンテキストのカーネルを線形結合した以下を用いた。

$$f(\mathbf{x}_*) = \mathbf{a}^T (e_{img} \mathbf{k}_{img,*} + e_{cxt} \mathbf{k}_{cxt,*}) + b$$

ここで、 e_m は m 番目のカーネルの重みとなり、 $\mathbf{k}_{m,*} = [k_m(\mathbf{x}_1, \mathbf{x}_*) \dots k_m(\mathbf{x}_N, \mathbf{x}_*)]^T$ である。ここで、 $m \in \{img, cxt\}$ であり、img と cxt はそれぞれ、画像とコンテキストを示す。画像分類に強い多項式カーネルを $\mathbf{k}_{img,*}$ として利用し、事前知識がないデータの分類によく使われる Radial Basis Function (RBF) を $\mathbf{k}_{cxt,*}$ として利用する。モデルのハイパーパラメータの設定に関しては文献 [13] を参照されたい。上述の識別関数を用いて、家電クラスごとに one-vs-rest SVM を用意し、テストインスタンスが分類されたクラスの中でマージンが最大となるクラスを最終的な推定結果、すなわち操作したい家電として出力する。また、MKL におけるパラメータの推定は Bayesian efficient multiple kernel learning (BE-MKL) を用いて行う。詳細は文献 [11] を参照されたい。

4.6.3 トレーニングデータが不十分な場合への対応

初期分類モデルの学習にはシステム導入段階で登録した家電画像のみを用いて行う。その後、日常生活の中で提案システムにより取得したラベリングされた画像データとラベルなしの行動・位置関連センサデータをトレーニングデータに追加し、分類モデルを更新する。そのため、システムの導入初期はトレーニングデータの量が不十分であるため、家電選択の精度が低くなる恐れがある。

この問題に対応するため、下記2つのアプローチを提案

する。

- 他のユーザのトレーニングデータの再利用：同じ環境で他のユーザが収集したラベルありデータを利用し、モデルの初期学習を行う。
- オンライン画像データベースの利用：インターネット上にある家電画像を活用することで問題に対応する。初期モデルの学習データには環境内の家電の画像に加えて、公開されている画像データベース内の家電画像も利用する。オンライン上の画像データベースには大量のラベル付きデータ (“television” など) を提供する ImageNet [7] を使用する。環境内にある家電に対応する画像はラベルを用いて ImageNet から検索・取得する。そして、モデルの学習には環境内の家電と類似する画像のみを抽出して利用する。類似画像の抽出には図2に示した構造のDCNNを用いる。DCNN から抽出した画像特徴量ベクトルのユークリッド距離を類似度として計算し、Top-k の類似オンライン画像のみを抽出する。

5. 評価実験

5.1 データセット

データの収集は、研究用途に建築されたスマートホーム内で行う。部屋内にある多くの電化製品や家具はホームネットワークを通じて操作できる。図4は実験環境のフロアマップを示す。表1は実験環境で利用した13種類の家電を示す、図4に各家電の環境内の位置も示す。3人の実験参加者は Glass とズボンポケットに入れたスマートフォンを用いて、実験環境でセンサデータを収集した。

本研究の提案手法は、ユーザの日常生活の中で利用されることを想定している。しかし、実際に実験参加者の家で長時間のカメラ画像付きデータを収集するのは困難であ

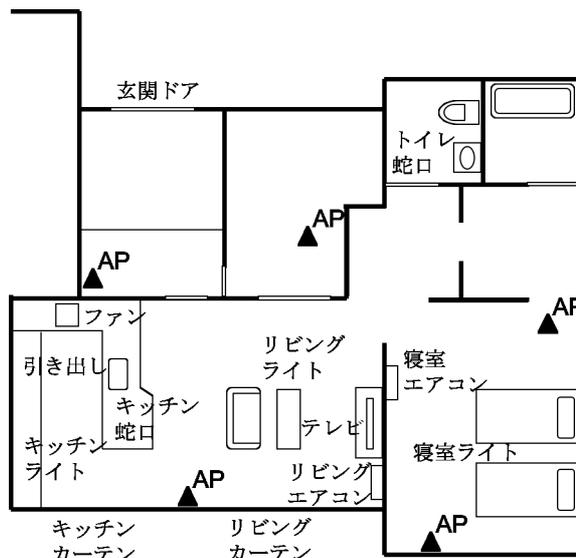


図4 実験環境のフロアマップ
Fig. 4 Floor plan of experimental environment.

表 1 実験で用いた家電

Table 1 Appliances installed in experimental environment.

リビングエアコン	寝室エアコン	寝室ライト
玄関ドア	引き出し	キッチン蛇口
換気扇	キッチンライト	リビングカーテン
キッチンカーテン	リビングライト	トイレ蛇口
テレビ		

る。したがって、本研究では semi-naturalistic collection protocol [1] を使用してセンサデータを収集する。Semi-naturalistic collection protocol では研究室でのデータ収集に比べて、実験参加者はある程度自由に行動できるため、多様なデータを収集できる。このプロトコルでは、参加者が指示書に従い、ランダムな順に行動を行う。指示書に書かれた指示は比較的あいまいであるため（たとえば、「TVを見る」、「トイレに行く」など）、実験参加者はそれぞれの行動を自由に行うことができる。また正解データを取得するため、参加者は頭をその方向に向け、3章で説明したアプリケーションを利用し、手動で該当家電の選択および操作を行う。Glass が記録した操作した家電の名前およびタイムスタンプからなるラベルは正解データとして用いる。

指示書には 8 種類の行動：「帰宅する」、「料理を準備する」、「食べる」、「食器を洗う」、「テレビを見る」、「トイレに行く」、「寝る」、「外出する」が記述されており、各行動中に家電の操作を行う。たとえば、料理を準備する際、「キッチンライト」、「ファン」、「引き出し」および「キッチン蛇口」などの家電を操作する。1 セッションのデータ収集は「帰宅する」から始まり、そして、「外出する」以外の 6 種類の行動をランダムな順で行い、最後に外出することで 1 セッションのデータ収集が終了する。各参加者は実験環境で 10 セッションのデータを収集した。また、家電の登録作業は実験を開始する前に行った。

5.2 評価手法

13 種類の家電を実験に用いたため、13 クラス分類問題となる。分類性能の評価には適合率、再現率および F 値 $(\frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}})$ を用いる。提案手法の効果を検証するため、評価実験では以下の 6 つの手法を用意・比較する。

- SVM w/ cam: 単純に画像の特徴のみを SVM に入力する (実装には LIBSVM [2] を用いた)。
- SVM all: 画像の特徴に加えて、コンテキスト情報の特徴も SVM に入力する。
- Proposed only cam: 提案手法。ただし、入力画像の特徴のみである。
- Proposed: 提案手法。
- Proposed w/o act: 提案手法。ただし、行動に関連するセンサデータは用いない。
- Proposed w/o pos: 提案手法。ただし、位置に関連する

表 2 評価手法の分類精度

Table 2 Classification accuracies for evaluation methods.

	precision	recall	F-measure
SVM all	76.6	70.0	73.2
SVM w/ cam	74.5	71.3	72.9
Proposed	85.8	78.1	81.0
Proposed only cam	83.5	72.0	75.1
Proposed w/o act	85.6	74.9	79.1
Proposed w/o pos	84.3	74.9	78.1
Proposed w/o IGMM	84.0	73.8	78.6

センサデータは用いない。

比較対象となる SVM は RBF カーネルを用いる。そのハイパーパラメータは LIBSVM [2] のデフォルト値を使用する。まず、家電登録フェーズの画像のみを用いて各手法の初期分類モデルを訓練した。また、各セッションの終了後、セッション中にラベリングされたデータをトレーニングデータに追加し、分類モデルを逐次的に更新する。これにより、ユーザがアプリケーションを用いて家電操作を繰り返すほど、家電選択の精度は向上する。

5.3 注視検出の結果

分類性能の評価結果を説明する前に、まず、本研究の注視検出方法の結果を説明する。収集したデータから注視を検出した直後、実験参加者が実際に家電の操作を行った場合、検出した注視は正しいと判断する。また、注視を検出した後、家電の操作を行うことなく次の注視が検出された場合、前者の注視は不正解と判断する。実際に実験で収集したデータセットを用いて計算した適合率と再現率はそれぞれ 70.4% と 94.6% となり、提案手法は高い再現率を達成した。すなわち、注視をほぼもれなく検出することができた。

5.4 家電選択の結果

・分類性能の評価

各セッション終了後、セッション中にラベリングされたデータをトレーニングデータに追加してモデルのアップデートを行う。表 2 は各手法における全セッションの平均適合率、再現率、F 値を示す。まず Proposed only cam に着目すると、画像のみを用いるだけでも良好な分類精度を達成できており、平均 F 値は 75% を超えていた。

表 2 の Proposed w/o act と Proposed w/o pos の結果に着目すると、位置情報が分類精度の向上に大きく貢献していることが分かる。位置情報を用いることで、「ライト」と「エアコン」の F 値が Proposed w/o act よりおよそ 5% 向上していた。また行動情報を用いると、「キッチンライト」と「ファン」などの精度が約 10% 向上した。「料理を準備する」ときに計測された行動に関するセンサデータ（音声、加速度）がこれらの家電の認識に有用であったと考える。

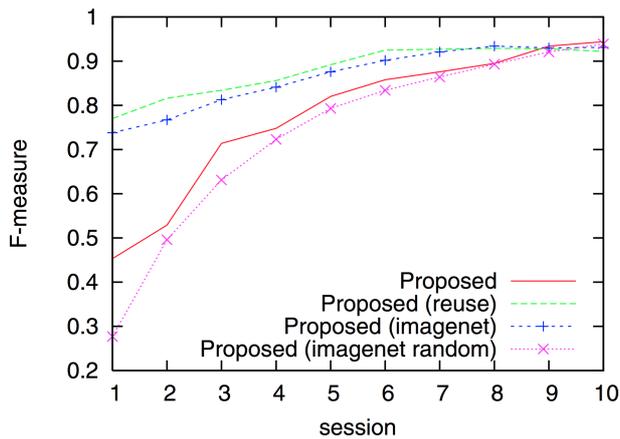


図 5 Proposed のテストセッションと平均 F 値の関係

Fig. 5 Transitions of average F-measures for Proposed method.

表 3 他のユーザのトレーニングデータの再利用による分類精度

Table 3 Classification results when we reuse other users' training data.

	precision	recall	F-measure
SVM all	82.7	81.5	82.1
SVM w/ cam	81.0	80.3	80.6
Proposed	93.5	90.2	91.8
Proposed only cam	90.8	86.4	88.5
Proposed w/o act	92.4	88.1	89.2
Proposed w/o pos	92.7	88.5	89.6

また、行動と位置情報を両方使った Proposed は 81.8%の精度を達成した。表 2 は IGMM を用いない場合 (Proposed w/o IGMM), すなわち抽出した各特徴量をそのまま MKL の学習に用いたときの結果も示した。その結果から IGMM の有効性を確認した (3.2%の精度向上)。

・トレーニングデータの量の影響

図 5 は提案手法におけるそれぞれのセッションの平均 F 値の遷移である。セッション 1 では、提案手法は導入時に撮影された画像のみを用いて、初期モデルを学習している。F 値はトレーニングデータの量の増加にともない上昇し、最終的には 95%の精度を達成した。

・他のユーザのトレーニングデータの再利用

上述のように、トレーニングデータの累積量が少ないセッションに関して、F 値が高いとはいえなかった。この問題に対処するため、本研究では同じ環境で他のユーザが収集したラベリングデータを再利用する。ここでは、1 人の実験参加者が収集したトレーニングデータに加えて、他の 2 人の参加者のトレーニングデータ (合計 20 セッション) を利用してモデルの学習を行う。表 3 にその結果を示す。また、図 5 の Proposed (reuse) は他の参加者のデータを再利用した場合の平均 F 値の遷移を示す。セッション 1 に着目すると、他のユーザの 20 セッションのデータを使用して初期モデルを学習することにより、セッション 1 の



図 6 導入段階撮影した家電画像と ImageNet から取得した類似画像およびその非類似度

Fig. 6 Photos taken during the installation period and their similar images obtained from ImageNet.

F 値を約 78%にまで向上させることができた。

・オンライン画像データベースの利用効果

ImageNet から取得したオンライン画像を利用した効果を紹介する。環境内にある各種類の家電に対応するオンライン画像を ImageNet から取得し、Top-k の類似画像を抽出する (k = 100)。抽出された類似画像は、実験参加者から収集した学習データに加えて、家電選択モデルの学習に用いる。図 6 は初期導入段階で撮影した家電画像に対する top-4 類似画像、および類似画像と撮影した画像との距離を示す。図 5 の Proposed (imagenet) はオンライン画像を用いた場合のセッションの平均 F 値の遷移を示す。セッション 1 の F 値は 73.8%となり、Proposed に比べて約 30%を向上させることができた。また、図 5 にはランダムで選択したオンライン画像をトレーニングデータとして用いた場合の精度も示す (Proposed (imagenet random))。精度はむしろ低下しており、類似する画像を用いることの有効性が示された。上述のように、オンライン画像を利用することで、システムの導入時から、家電選択の精度を向上させることができる。

・Leave one session out 交差検定

トレーニングデータが十分にある場合の、提案手法の性能を検証する。評価手法には leave-one-session-out 交差検定を用いる。これは、1 セッションのデータを除いて、残る 9 セッションのデータから訓練したモデルを、除いた 1 セッションのデータでテストする交差検定手法である。表 4 にその結果を示す。9 セッションのデータを用いた場合、提案手法は約 95%の高い分類性能を示し、図 5 における 10 セッション目の精度とほぼ同じ精度を達成した。そして、行動および位置情報を利用することで、画像のみを利用した場合と比べて、精度がおよそ 10%改善された。

図 7 に Proposed only cam の混同行列の結果を示す。行

表 4 Leave one session out 交差検定による分類精度

Table 4 Classification accuracies using leave one session cross validation.

	precision	recall	F-measure
SVM all	84.5	84.4	84.4
SVM w/ cam	81.3	81.2	81.2
Proposed	95.5	93.6	94.5
Proposed only cam	85.7	86.2	85.9
Proposed w/o act	92.8	89.7	91.2
Proposed w/o pos	89.4	87.8	88.6

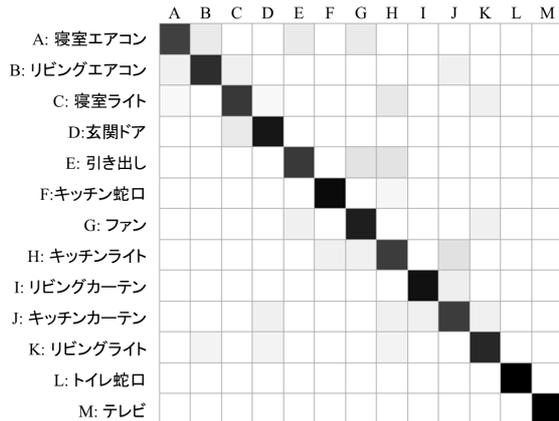


図 7 Proposed only cam の混同行列

Fig. 7 Visual confusion matrix of Proposed only cam result.

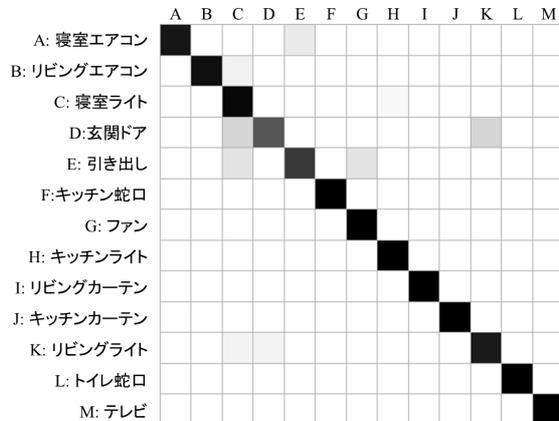


図 8 Proposed の混同行列

Fig. 8 Visual confusion matrix of Proposed result.

列に示すように、「エアコン」と「ライト」の分類結果は比較的悪く、Proposed only cam 手法は「キッチンライト」と「寝室ライト」をうまく区別できていなかった。Proposed only cam 手法はカメラ画像のみを用いるため、似た外見を持つ家電をうまく区別できないと考える。また、「引き出し」に関する精度も悪かった。これは、「引き出し」の外見が非常にシンプルで、区別するのに有用な画像の特徴が少ないためと考える。

図 8 に Proposed の混同行列の結果を示す。Proposed only cam における「エアコン」と「ライト」の結果と比べ

ると、精度はおおよそ 14% 向上した。また表 4 に示したように、コンテキスト情報を利用した Proposed は Proposed only cam に比べると、大幅に精度が改善され、約 10% 向上した。一方、SVM all は画像のみを利用する SVM w/ cam と比べると、F 値が約 3.5% 向上した。通常の SVM を用いた手法に比べて、異なるソースのデータの扱いに長けた MKL を用いた提案手法の性能が高いことも確認できる。また、Proposed only cam の性能は SVM w/ cam より高かった。これは、Proposed only cam が画像の扱いに適したカーネル関数を用いているためと考える。図 7 と比較すると、「玄関ドア」と「リビングライト」間の誤判定が増えた。これは、「玄関ドア」と「リビングライト」を選択するときの行動と位置情報が類似するためだと考える。具体的には、「リビングライト」を操作する場合、玄関から入ってきてすぐに操作する場合と、玄関から出ていく前に操作する場合があり、「玄関ドア」と近い位置で操作されており、操作する前後の行動もともに「歩く」行動である。このコンテキスト情報の類似性により、推定結果に影響が出たと考えられる。「引き出し」と「寝室ライト」間の誤判定が増えた原因は、実験参加者らが寝室以外のリビングの寝室ライトが見える場所から「寝室ライト」を選択・操作を行ったため、位置情報の学習がうまくできなかったからと考えられる。

・計算時間

提案手法の処理時間について紹介する。特徴量の抽出は並列で行い、そして画像特徴量の計算はサーバで行うため、全体の処理時間は画像特徴量の計算時間に依存する。画像の転送および特徴量計算時間はそれぞれ 0.33 秒と 0.08 秒である。また、分類の計算時間は 0.01 秒であるため、全体の処理時間はおおよそ 0.43 秒となる。

5.5 IGMM の結果

ノンパラメトリッククラスタリング手法を用いて行動および位置データをクラスタリングする IGMM の結果を紹介する。IGMM は位置推定および行動推定用にそれぞれ用意しており、それぞれに対応するセンサデータを類似度に基づきクラスタ化し、ある時刻（たとえば家電への注視が発生する時刻）のセンサデータがどのクラスタに属するかを推定することができる。たとえば行動の IGMM では、「調理」や「睡眠」などに関わるクラスタがセンサデータから自動的に形成され、ある時刻のユーザのデータがどのクラスタに属するかが分かる。また位置の IGMM では、「キッチン」や「ベッドルーム」といった位置に関するクラスタが形成される。図 9 と図 10 は、被験者が様々な日常行動を行ったときに得られたセンサデータを IGMM によりクラスタリングし、主成分分析 (PCA) を用いて、2次元のデータに次元圧縮してプロットした結果である。図中の各ポイントは実験で取得したラベルなしのセンサデータポイ

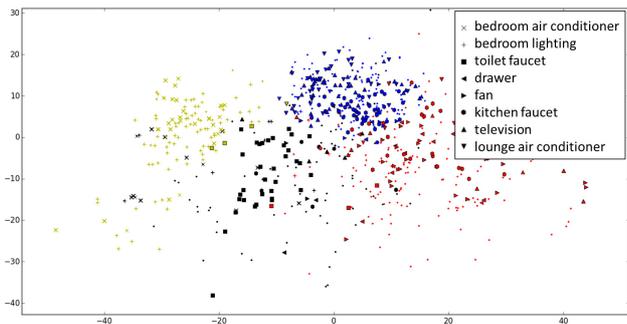


図 9 位置推定用 IGMM の結果

Fig. 9 Clustering result of IGMM for positional information.

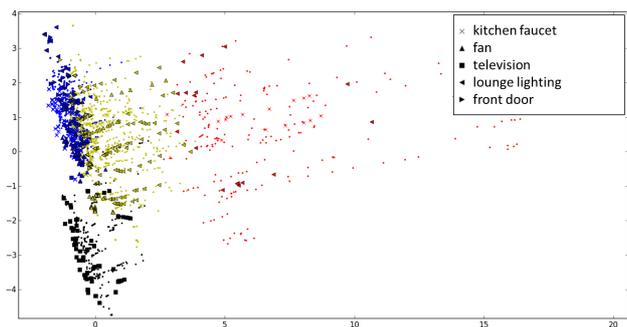


図 10 行動推定用 IGMM の結果

Fig. 10 Clustering result of IGMM for activity information.

ントとなる。ポイントの色はそのデータがどのクラスに属するかを示す。また、図 9 と図 10 には、代表的な家電を操作したときに得られたデータポイントも示す。図 9 に示したように、同じ場所で使う家電に関連する多くのデータポイントは同じクラスに分類される。たとえば、図 9 のベッドルームにある家電を操作した際に得られたコンテキストデータが同じクラス（黄色のクラス）に分類されていることが分かる。図 10 でも、「調理」に関する家電を操作した際に得られたコンテキストデータが同じクラス（青色のクラス）に分類されていることが分かる。また、「テレビを見る」など簡単な動きの行動に対応するデータポイントは、複雑な動きをとる「調理」の行動に対応するデータポイントとの距離が大きいことが分かる。

6. 考察

本論文では、スマートグラスを用いた家電操作のための家電選択手法の提案およびそれを実現するセンサデータ処理手法を提案した。提案した選択手法では、たとえば既存の万能リモコンなどと異なり、ハンズフリーな家電選択を可能とするため、家事中やリモコンが手元にない場合に有効である。さらに Amazon Echo などの設置型デバイスと異なりユーザの場所を選ばない。本論文で実現したスマートグラスを用いた家電選択方法について、そのユーザビリティを評価するため、今後下記の検証を行うことが必要であると考えられる。(1) 頭部の向きから家電を選択するアプ

ローチを用いているため、実際の生活の中で向きの変更が困難な状況およびその状況への対応方法を検証する。(2) 加速度などのセンサデータの常時収集が必要なため、その常時収集によるデバイス稼働時間およびユーザ体験への影響を確認する。(3) 本提案手法を既存手法（リモコン、音声など）と比較し、それぞれ利点・欠点を実生活におけるユーザによる長時間的な利用から検証する。

7. まとめ

本論文では Google Glass を用いて行動や位置といったコンテキスト情報に基づく新しい家電選択手法とその実現方法を提案した。提案手法の有効性を検証するための評価実験は実際のスマートホームで行い、トレーニングデータが十分にある場合は平均 91.8% の精度で家電を正しく選択できた。

提案した手法には、以下の制約が考えられる。(1) 同時利用の制限：同じ部屋にいる複数のユーザが同時に同じ家電を操作しようとした際にコンフリクトが発生する。この問題に関しては、ユーザ数の同時アクセス制限などの機能を追加することで解決できると考える。(2) ポータブル家電の認識：提案手法は家電の位置情報を用いて認識を行うため、ポータブル家電の使用位置が大きく変化する場合に、認識精度に影響する恐れがある。(3) 利用頻度が少ない家電の認識：学習データはユーザの日常生活の中で蓄積されるため、使用頻度が少ない家電に対して、学習データが十分に集まらない恐れがある。学習におけるクラス不均衡問題への対応が必要になる。(4) 複数家電が近くに配置される場合、ユーザが最も選択したい家電への認識精度が低下する傾向がある。この問題について、本提案手法は画像情報に加えてコンテキスト情報により最も選択したい家電の推定性能を向上させることが可能と考える。また、複数の家電がそれぞれ認識された場合、家電の同時複数選択を提示し、それらの連携操作を行わせることも可能と考える。たとえば、録画装置とテレビが認識された場合、「録画した動画をテレビ上で見る」操作がユーザに提示される。上記家電連携操作方法の実現と提案手法のユーザビリティに関する考察は今後の重要な課題の 1 つである。

謝辞 本研究は JST CREST JPMJCR15E2, JSPS 科研費 JP26730047 の助成を受けたものであり、総務省委託研究「脳の仕組みを活かしたイノベーション創成型研究開発」, 革新的研究開発推進プログラム (ImPACT) の一環として実施したものです。

参考文献

- [1] Bao, L. and Intille, S.S.: Activity recognition from user-annotated acceleration data, *Pervasive 2004*, pp.1-17 (2004).
- [2] Chang, C.-C. and Lin, C.-J.: LIBSVM: A library for support vector machines, *ACM Trans. Intelligent Systems*

and Technology, Vol.2, pp.27:1–27:27 (2011).

[3] Chen, T., Morris, J. and Martin, E.: Probability density estimation via an infinite Gaussian mixture model: Application to statistical process monitoring, *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, Vol.55, No.5, pp.699–715 (2006).

[4] Christensen, H., Casanueva, I.N., Cunningham, S., Green, P. and Hain, T.: homeService: Voice-enabled assistive technology in the home using cloud-based automatic speech recognition, *4th Workshop on Speech and Language Processing for Assistive Technologies*, pp.29–34 (2013).

[5] Corno, F., Gale, A., Majaranta, P. and Rähkä, K.-J.: Eye-based Direct Interaction for Environmental Control in Heterogeneous Smart Environments, *Handbook of Ambient Intelligence and Smart Environments*, pp.1117–1138 (2010).

[6] Cowling, M.: Non-speech environmental sound recognition system for autonomous surveillance, Ph.D. Thesis, Griffith University (2004).

[7] Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K. and Fei-Fei, L.: ImageNet: A Large-Scale Hierarchical Image Database, *CVPR 2009*, pp.248–255 (2009).

[8] Donahue, J., Jia, Y., Vinyals, O., Hoffman, J., Zhang, N., Tzeng, E. and Darrell, T.: DeCAF: A deep convolutional activation feature for generic visual recognition, arXiv preprint arXiv:1310.1531 (2013).

[9] Ferguson, T.S.: A Bayesian analysis of some nonparametric problems, *The Annals of Statistics*, pp.209–230 (1973).

[10] Gafurov, D., Helkala, K. and Sondrol, T.: Biometric gait authentication using accelerometer sensor, *Journal of computers*, Vol.1, No.7, pp.51–59 (2006).

[11] Gönen, M.: Bayesian Efficient Multiple Kernel Learning, *ICML 2012* (2012).

[12] Hardegger, M., Tröster, G. and Roggen, D.: Improved ActionSLAM for long-term indoor tracking with wearable motion sensors, *International Symposium on Wearable Computers (ISWC2013)*, pp.1–8 (2013).

[13] Hastie, T., Tibshirani, R. and Friedman, J.: Kernel Smoothing Methods, *The Elements of Statistical Learning*, Springer, pp.191–218 (2009).

[14] Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S. and Darrell, T.: Caffe: Convolutional architecture for fast feature embedding, *ACM Multimedia 2014*, pp.675–678 (2014).

[15] Lane, N.D., Xu, Y., Lu, H., Hu, S., Choudhury, T., Campbell, A.T. and Zhao, F.: Enabling large-scale human activity inference on smartphones using community similarity networks (CSN), *UbiComp 2011*, pp.355–364 (2011).

[16] Maekawa, T. and Watanabe, S.: Unsupervised Activity Recognition with User’s Physical Characteristics Data, *International Symposium on Wearable Computers (ISWC 2011)*, pp.89–96 (2011).

[17] Mohan, A., Woo, G., Hiura, S., Smithwick, Q. and Raskar, R.: Bokode: Imperceptible Visual Tags for Camera Based Interaction from a Distance, *ACM Trans. Graph.*, Vol.28, No.3, pp.98:1–98:8 (2009).

[18] Moran, T.P., Saund, E., van Melle, W., Gujar, A., Fishkin, K.P. and Harrison, B.L.: Design and Technology for Collaborative: Collaborative Collages of Information on Physical Walls, *ACM Symposium on User Interface Software and Technology (UIST 1999)*, pp.197–206 (1999).

[19] Neßelrath, R., Lu, C., Schulz, C.H., Frey, J. and Alexandersson, J.: A Gesture Based System for Context-Sensitive Interaction with Smart Homes, *Ambient Assisted Living Congress 2011*, pp.209–219 (2011).

[20] Rekimoto, J. and Ayatsuka, Y.: CyberCode: Designing Augmented Reality Environments with Visual Tags, *Designing Augmented Reality Environments (DARE 2000)*, pp.1–10 (2000).

[21] Shi, F., Gale, A. and Purdy, K.: Helping People with ICT Device Control by Eye Gaze, *Computers Helping People with Special Needs*, pp.480–487 (2006).

[22] Solanki, U.V. and Desai, N.H.: Hand gesture based remote control for home appliances: Handmote, *World Congress on Information and Communication Technologies (WICT 2011)*, pp.419–423 (2011).

[23] Tsukada, K. and Yasumura, M.: Ubi-finger: Gesture input device for mobile use, *Ubicomp 2001*, p.11 (2001).



孔全 (正会員)

2014年3月大阪大学大学院情報科学研究科博士前期課程修了。2016年3月同大学院情報科学研究科博士後期課程修了。同年より株式会社日立製作所研究開発グループに入社。現在メディア情報処理に関する研究に従事。

ACM 会員。



前川卓也 (正会員)

2003年大阪大学工学部電子情報エネルギー工学科卒業。2006年同大学院情報科学研究科博士後期課程修了。同年日本電信電話株式会社入社。2012年4月より大阪大学大学院情報科学研究科准教授。2013年8～10月スイス連邦工科大学ローザンヌ校招聘教授。博士(情報科学)。本会平成22年度山下記念研究賞、日本データベース学会平成25年度上林奨励賞等受賞。ACM, IEEE, 電気学会, 日本データベース学会各会員。



宮西大樹 (正会員)

2011年神戸大学大学院工学研究科情報知能学専攻博士前期課程修了。2014年同大学院システム情報学研究科計算科学専攻博士後期課程修了。現在、国際電気通信基礎技術研究所研究員。情報検索、知能情報システムの研究に従事。



須山 敬之 (正会員)

1992年大阪大学大学院機械工学専攻博士前期課程修了。NTT コミュニケーション科学研究所，NTT 西日本研究開発センタ等を経て，2014年より国際電気通信基礎技術研究所に所属。現在，動的脳イメージング研究室長。博士（情報学）。センサネットワーク，ユビキタスコンピューティング，ブレイン・マシン・インタフェースの研究に従事。電子情報通信学会，IEEE，ACM 各会員。