

(12) 特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関
国際事務局

(43) 国際公開日
2020年12月3日(03.12.2020)



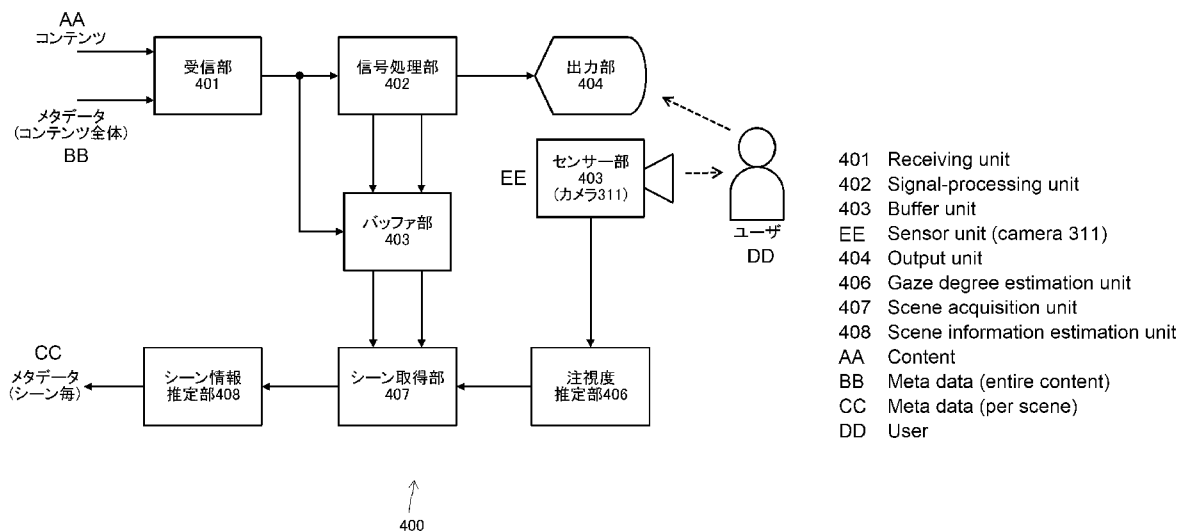
(10) 国際公開番号
WO 2020/240976 A1

- (51) 国際特許分類:
G06F 16/9035 (2019.01) H04N 21/442 (2011.01)
G06F 16/9038 (2019.01) H04N 21/466 (2011.01)
G06F 16/907 (2019.01)
- (21) 国際出願番号: PCT/JP2020/009957
- (22) 国際出願日: 2020年3月9日(09.03.2020)
- (25) 国際出願の言語: 日本語
- (26) 国際公開の言語: 日本語
- (30) 優先権データ:
特願 2019-098472 2019年5月27日(27.05.2019) JP
- (71) 出願人: ソニー株式会社 (SONY CORPORATION) [JP/JP]; 〒1080075 東京都港区港南1丁目7番1号 Tokyo (JP).
- (72) 発明者: 高木 悟郎 (TAKAKI, Goro); 〒1080075 東京都港区港南1丁目7番1号 ソニー株式会社内 Tokyo (JP). 小林 由幸 (KOBAYASHI, Yoshiyuki); 〒1080075 東京都港区港南1丁目7番1号 ソニー株式会社内 Tokyo (JP).
- (74) 代理人: 宮田 正昭, 外 (MIYATA, Masaaki et al.); 〒1040032 東京都中央区八丁堀三丁目25番9号 Daiwa八丁堀駅前ビル西館8階 特許業務法人 大同特許事務所 Tokyo (JP).
- (81) 指定国(表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JO, JP, KE, KG, KH,

(54) Title: ARTIFICIAL INTELLIGENCE INFORMATION PROCESSING DEVICE AND ARTIFICIAL INTELLIGENCE INFORMATION PROCESSING METHOD

(54) 発明の名称: 人工知能情報処理装置及び人工知能情報処理方法

[図4]



(57) Abstract: The present invention provides an artificial intelligence information processing device that generates information related to scenes by artificial intelligence. The artificial intelligence information processing device comprises a gaze degree estimation unit for estimating the gaze degree of a user viewing content by artificial intelligence on the basis of sensor information, an acquisition unit for acquiring the video of a scene in the content which the user is gazing at and information relating to the content on the basis of the estimation result of the gaze degree estimation unit, and a scene information estimation unit for estimating information relating to the scene which the user is gazing at by artificial intelligence on the basis of the video of the scene which the user is gazing at and information relating to the content.

WO 2020/240976 A1

KN, KP, KR, KW, KZ, LA, LC, LK, LR, LS, LU, LY,
MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ,
NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT,
QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL,
ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG,
US, UZ, VC, VN, WS, ZA, ZM, ZW.

- (84) 指定国(表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ, BY, KG, KZ, RU, TJ, TM), ヨーロッパ (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

添付公開書類：

- 一 国際調査報告 (条約第21条(3))

(57) 要約：シーンに関する情報を人工知能により生成する人工知能情報処理装置を提供する。人工知能情報処理装置は、コンテンツを視聴しているユーザの注視度をセンサー情報に基づいて人工知能により推定する注視度推定部、前記注視度推定部の推定結果に基づいて、前記コンテンツ中で前記ユーザが注視するシーンの映像と前記コンテンツに関する情報を取得する取得部、前記ユーザが注視するシーンの映像と前記コンテンツに関する情報に基づいて前記ユーザが注視しているシーンに関する情報を人工知能により推定するシーン情報推定部を具備する。

明 細 書

発明の名称：人工知能情報処理装置及び人工知能情報処理方法

技術分野

[0001] 本明細書で開示する技術は、コンテンツに関する情報を人工知能により処理する人工知能情報処理装置及び人工知能情報処理方法に関する。

背景技術

[0002] テレビ放送サービスが広範に普及して久しい。現在、テレビ受信機は広範に普及しており、各家庭に1台又は複数台設置されている。最近では、IPTV (Internet Protocol TV) やOTT (Over-The-Top) といった、ネットワークを利用した放送型の動画配信サービスも浸透しつつある。

[0003] また最近では、テレビ受信機とセンシング技術とを組み合わせ、視聴者の映像コンテンツに対する注視度合いを示す「視聴質」を計測する技術についても研究開発がなされている（例えば、特許文献1を参照のこと）。視聴質の利用方法はさまざまである。例えば、視聴質の計測結果に基づいて、映像コンテンツや広告の効果を評価したり、視聴者に他のコンテンツや商品の推薦を行ったりすることができる。

先行技術文献

特許文献

[0004] 特許文献1：WO2017/120469

特許文献2：特許第4840393号公報

特許文献3：特開2007-143010号公報

特許文献4：特開2008-236779号公報

発明の概要

発明が解決しようとする課題

[0005] 本明細書で開示する技術の目的は、コンテンツに付随する情報を人工知能により処理する人工知能情報処理装置及び人工知能情報処理方法を提供する

ことにある。

課題を解決するための手段

- [0006] 本明細書で開示する技術の第1の側面は、
コンテンツを視聴しているユーザの注視度をセンサー情報に基づいて人工知能により推定する注視度推定部と、
前記注視度推定部の推定結果に基づいて、前記コンテンツ中で前記ユーザが注視するシーンの映像と前記コンテンツに関する情報を取得する取得部と、
、
前記ユーザが注視するシーンの映像と前記コンテンツに関する情報に基づいて、前記ユーザが注視しているシーンに関する情報を人工知能により推定するシーン情報推定部と、
を具備する人工知能情報処理装置である。
- [0007] 前記シーン情報推定部は、人工知能による推定として、シーンの映像と前記コンテンツに関する情報と、シーンに関する情報との相関関係を学習したニューラルネットワークを利用して、前記ユーザが注視するシーンと相関関係のある情報を推定する。
- [0008] また、前記注視度推定部は、人工知能による推定として、センサー情報と前記ユーザの注視度との相関関係を学習したニューラルネットワークを利用して、前記コンテンツを視聴している前記ユーザに関するセンサー情報と相関関係のある注視度合いを推定する。
- [0009] また、本明細書で開示する技術の第2の側面は、
コンテンツを視聴しているユーザに関するセンサー情報を入力する入力部と、
センサー情報、コンテンツ及びコンテンツの情報と、ユーザが注視するシーンに関する情報との相関関係を学習したニューラルネットワークを利用して、前記ユーザが注視するシーンと相関関係のある情報を推定するシーン情報推定部と、
を具備する人工知能情報処理装置である。

- [0010] また、本明細書で開示する技術の第3の側面は、
コンテンツを視聴しているユーザの注視度をセンサー情報に基づいて人工
知能により推定する注視度推定ステップと、
前記注視度推定ステップにおける推定結果に基づいて、前記コンテンツ中
で前記ユーザが注視するシーンの映像と前記コンテンツに関する情報を取得
する取得ステップと、
前記ユーザが注視するシーンの映像と前記コンテンツに関する情報に基づ
いて、前記ユーザが注視しているシーンに関する情報を人工知能により推定
するシーン情報推定ステップと、
を有する人工知能情報処理方法である。

発明の効果

- [0011] 本明細書で開示する技術によれば、コンテンツの一部のシーンに関するメ
タデータを人工知能により推定する人工知能情報処理装置及び人工知能情報
処理方法を提供することができる。
- [0012] なお、本明細書に記載された効果は、あくまでも例示であり、本明細書で
開示する技術によりもたらされる効果はこれに限定されるものではない。ま
た、本明細書で開示する技術が、上記の効果以外に、さらに付加的な効果を
奏する場合もある。
- [0013] 本明細書で開示する技術のさらに他の目的、特徴や利点は、後述する実施
形態や添付する図面に基づくより詳細な説明によって明らかになるであらう
。

図面の簡単な説明

- [0014] [図1]図1は、映像コンテンツを視聴するシステムの構成例を示した図である
。
[図2]図2は、テレビ受信装置100の構成例を示した図である。
[図3]図3は、テレビ受信装置100に装備されるセンサー群300の構成例
を示した図である。
[図4]図4は、シーン取得及びシーン情報推定システム400の構成例を示し

た図である。

[図5]図5は、注視度推定部406で利用されるニューラルネットワーク500の構成例を示した図である。

[図6]図6は、シーン情報推定部408で利用されるニューラルネットワーク600の構成例を示した図である。

[図7]図7は、シーン取得及びシーン情報推定システム400で実施される処理手順を示したフローチャートである。

[図8]図8は、シーン取得及びシーン情報推定システム800の変形例を示した図である。

[図9]図9は、シーン情報推定部806で利用されるニューラルネットワーク900の構成例を示した図である。

[図10]図10は、テレビを視聴するユーザをカメラで撮影する様子を示した図である。

[図11]図11は、カメラの撮影画像から顔認識する様子を示した図である。

[図12]図12は、ユーザの注視度の推測結果に伴って注視シーンを捕捉する様子を示した図である。

[図13]図13は、ユーザが注視したシーンの一覧画面の構成例を示した図である。

[図14]図14は、シーンを選択した後に遷移する画面の構成例示した図である。

[図15]図15は、関連コンテンツを提示する画面の構成例を示した図である。

[図16]図16は、ユーザが注視したシーンをスマートフォンの画面上でフィードバックする様子を示した図である。

[図17]図17は、複数のシーンの配置例を示した図である。

[図18]図18は、複数のシーンの配置例を示した図である。

[図19]図19は、パネルスピーカ技術の適用例を示した図である。

[図20]図20は、クラウドを利用した人工知能システム2000の構成例を

示した図である。

発明を実施するための形態

[0015] 以下、図面を参照しながら本明細書で開示する技術の実施形態について詳細に説明する。

[0016] A. システム構成

図1には、映像コンテンツを視聴するシステムの構成例を模式的に示している。

[0017] テレビ受信装置100は、映像コンテンツを表示する大画面並びの音声を出力するスピーカーを装備している。テレビ受信装置100は、例えば放送信号を選局受信するチューナーを内蔵し、又はセットトップボックスが接続されており、テレビ局が提供する放送サービスを利用することができる。放送信号は、地上波及び衛星波のいずれを問わない。

[0018] また、テレビ受信装置100は、例えばIPTVやOTTといったネットワークを利用した放送型の動画配信サービスも利用することができる。このため、テレビ受信装置100は、ネットワークインターフェースカードを装備し、イーサネット（登録商標）やWi-Fi（登録商標）などの既存の通信規格に基づく通信を利用して、ルータ経由やアクセスポイント経由でインターネットなどの外部ネットワークに相互接続されている。

[0019] インターネット上には、映像ストリームを配信するストリーム配信サーバが設置されており、テレビ受信装置100に対して放送型の動画配信サービスを提供する。

[0020] また、インターネット上には、さまざまなサービスを提供する無数のサーバが設置されている。サーバの一例は、ストリーム配信サーバである。テレビ受信装置100側では、ブラウザ機能を起動し、ストリーム配信サーバに対して例えばHTTP（Hyper Text Transfer Protocol）リクエストを発行して、Webサービスを利用することができる。

[0021] また、本実施形態では、クライアントに対してインターネット上で（若し

くは、クラウド上で)人工知能の機能を提供する人工知能サーバ(図示しない)も存在することを想定している。ここで、人工知能の機能とは、例えば、学習、推論、データ創出、計画立案といった、一般的に人間の脳が発揮する機能をソフトウェア又はハードウェアによって人工的に実現した機能を指す。また、人工知能サーバは、例えば、人間の脳神経回路を模したモデルにより深層学習(Deep Learning:DL)を行うニューラルネットワークを搭載している。ニューラルネットワークは、シナプスの結合によりネットワークを形成した人工ニューロン(ノード)が、学習によりシナプスの結合強度を変化させながら、問題に対する解決能力を獲得する仕組みを備えている。ニューラルネットワークは、学習を重ねることで、問題に対する解決ルールを自動的に推論することができる。なお、本明細書で言う「人工知能サーバ」は、単一のサーバ装置とは限らず、例えばクラウドコンピューティングサービスを提供するクラウドの形態であってもよい。

[0022] 図2には、テレビ受信装置100の構成例を示している。テレビ受信装置100は、主制御部201と、バス202と、ストレージ部203と、通信インターフェース(IF)部204と、拡張インターフェース(IF)部205と、チューナー/復調部206と、デマルチプレクサ(DEMUX)207と、映像デコーダ208と、音声デコーダ209と、文字スーパーデコーダ210と、字幕デコーダ211と、字幕合成部212と、データデコーダ213と、キャッシュ部214と、アプリケーション(AP)制御部215と、ブラウザ部216と、音源部217と、映像合成部218と、表示部219と、音声合成部220と、音声出力部221と、操作入力部222を備えている。

[0023] 主制御部201は、例えばコントローラとROM(Read Only Memory)(但し、EEPROM(Electrically Erasable Programmable ROM)のような書き換え可能なROMを含むものとする)、及びRAM(Random Access Memory)で構成され、所定の動作プログラムに従ってテレビ受信装置1

00全体の動作を統括的に制御する。コントローラは、CPU (Central Processing Unit)、MPU (Micro Processing Unit)、又はGPU (Graphics Processing Unit) 若しくはGPGPU (General Purpose Graphic Processing Unit) などROMは、オペレーティングシステム (OS) などの基本動作プログラムやその他の動作プログラムが格納された不揮発性メモリである。ROM内には、テレビ受信装置100の動作に必要な動作設定値が記憶されてもよい。RAMはOSやその他の動作プログラム実行時のワークエリアとなる。バス202は、主制御部201とテレビ受信装置100内の各部との間でデータ送受信を行うためのデータ通信路である。

[0024] ストレージ部203は、フラッシュROMやSSD (Solid State Drive)、HDD (Hard Disc Drive) などの不揮発性の記憶デバイスで構成される。ストレージ部203は、テレビ受信装置100の動作プログラムや動作設定値、テレビ受信装置100を使用するユーザの個人情報などを記憶する。また、インターネットを介してダウンロードした動作プログラムやその動作プログラムで作成した各種データなどを記憶する。また、ストレージ部203は、放送波やインターネットを通じて取得した動画、静止画、音声などのコンテンツも記憶可能である。

[0025] 通信インターフェース部204は、ルータ (前述) を介してインターネットと接続され、インターネット上の各サーバ装置やその他の通信機器とデータの送受信を行う。また、通信回線を介して伝送される番組のデータストリームの取得も行うものとする。ルータとは、イーサネット (登録商標) などの有線接続、あるいはWi-Fi (登録商標) などの無線接続のいずれであってもよい。

[0026] チューナー／復調部206は、アンテナ (図示しない) を介して地上波放送又は衛星放送などの放送波を受信し、主制御部201の制御に基づいてユーザの所望するサービス (放送局など) のチャンネルに同調 (選局) する。

また、チューナー／復調部 206 は、受信した放送信号を復調して放送データストリームを取得する。なお、複数画面同時表示や裏番組録画などを目的として、テレビ受信装置 100 が複数のチューナー／復調部を搭載する構成（すなわち多重チューナ）であってもよい。

[0027] デマルチプレクサ 207 は、入力したデータストリーム中の制御信号に基づいてリアルタイム提示要素である映像データストリーム、音声データストリーム、文字スーパーデータストリーム、字幕データストリームを、それぞれ映像デコーダ 208、音声デコーダ 209、文字スーパーデコーダ 210、字幕デコーダ 211 に分配する。デマルチプレクサ 207 に入力されるデータストリームは、放送サービスによる放送データストリームや、IPTV や OTT などの配信サービスによる配信データストリームを含む。前者は、チューナー／復調部 206 で選局受信及び復調された後にデマルチプレクサ 207 に入力され、後者は、通信インターフェース部 204 で受信された後にデマルチプレクサ 207 に入力される。また、デマルチプレクサ 207 は、マルチメディアアプリケーションやその構成要素であるファイル系データを再生し、アプリケーション制御部 215 に出力し、又はキャッシュ部 214 で一時的に蓄積する。

[0028] 映像デコーダ 208 は、デマルチプレクサ 207 から入力した映像ストリームを復号して映像情報を出力する。また、音声デコーダ 209 は、デマルチプレクサ 207 から入力した音声ストリームを復号して音声情報を出力する。デジタル放送では、例えば MPEG2 System 規格に則ってそれぞれ符号化された映像ストリーム並びに音声ストリームが多重化して伝送又は配信されている。映像デコーダ 208 並びに音声デコーダ 209 は、デマルチプレクサ 207 でデマルチプレクスされた符号化映像ストリーム、符号化映像ストリームを、それぞれ規格化されたデコード方式に従ってデコード処理を実施することになる。なお、複数種類の映像データストリーム及び音声データストリームを同時に復号処理するために、テレビ受信装置 100 は複数の映像デコーダ 208 及び音声デコーダ 143 を備えてもよい。

- [0029] 文字スーパーデコーダ210は、デマルチプレクサ207から入力した文字スーパーデータストリームを復号して文字スーパー情報を出力する。字幕デコーダ211は、デマルチプレクサ207から入力した字幕データストリームを復号して字幕情報を出力する。字幕合成部212は、文字スーパーデコーダ210から出力された文字スーパー情報と、字幕デコーダ211から出力された字幕情報は、字幕合成部212とを合成処理する。
- [0030] データデコーダ213は、MPEG-2 TSストリームに映像及び音声とともに多重化されるデータストリームをデコードする。例えば、データデコーダ213は、PSI (Program Specific Information) テーブルの1つであるPMT (Program Map Table) の記述子領域に格納された汎用イベントメッセージをデコードした結果を、主制御部201に通知する。
- [0031] アプリケーション制御部215は、放送データストリームに含まれる制御情報をデマルチプレクサ207から入力し、または、通信インターフェース部204を介してインターネット上のサーバ装置から取得して、これら制御情報を解釈する。
- [0032] ブラウザ部216は、キャッシュ部214若しくは通信インターフェース部204を介してインターネット上のサーバ装置から取得したマルチメディアアプリケーションファイルやその構成要素であるファイル系データを、アプリケーション制御部215の指示に従って提示する。ここで言うマルチメディアアプリケーションファイルは、例えばHTML (Hyper Text Markup Language) 文書やBML (Broadcast Markup Language) 文書などである。また、ブラウザ部216は、音源部217に働きかけることにより、アプリケーションの音声情報の再生も行うものとする。
- [0033] 映像合成部218は、映像デコーダ208から出力された映像情報と、字幕合成部212から出力された字幕情報と、ブラウザ部216から出力されたアプリケーション情報を入力し、適宜選択し又は重畳する処理を行う。映

像合成部 218 はビデオ RAM（図示を省略）を備え、このビデオ RAM に入力された映像情報に基づいて表示部 219 の表示駆動が実施される。また、映像合成部 218 は、主制御部 201 の制御に基づいて、必要に応じて、EPG（Electronic Program Guide）画面や、主制御部 201 が実行するアプリケーションによって生成されたグラフィックスなどの画面情報の重畳処理も行う。

[0034] 表示部 219 は、例えば液晶ディスプレイや有機 EL（Electro-Luminescence）ディスプレイなどからなる表示デバイスであり、映像合成部 218 で選択又は重畳処理を施された映像情報をユーザに提示する。また、表示部 219 として、透過型の液晶パネルを複数の表示領域（ブロック）に分割して、表示領域毎のバックライトを用いて個別に光を入射させるタイプの液晶表示装置を利用してもよい（例えば、特許文献 2 を参照のこと）。この種の表示装置によれば、表示領域毎に入射光の量を制御して表示画像の輝度のダイナミックレンジ拡大を実現するといった利点がある。

[0035] 音声合成部 220 は、音声デコーダ 209 から出力された音声情報と、音源部 217 で再生されたアプリケーションの音声情報を入力して、適宜選択又は合成などの処理を行う。

[0036] 音声出力部 221 は、チューナー／復調部 206 で選局受信した番組コンテンツやデータ放送コンテンツの音声出力や、音声合成部 220 で処理された音声情報（音声ガイダンス又は音声エージェントの合成音声などを含む）の出力に用いられる。音声出力部 221 は、スピーカーなどの音響発生素子で構成される。例えば、音声出力部 221 は、複数のスピーカーを組み合わせたスピーカーアレイ（多チャンネルスピーカー若しくは超多チャンネルスピーカー）であってもよく、一部又は全部のスピーカーがテレビ受信装置 100 に外付け接続されていてもよい。

[0037] コーン型スピーカーの他、フラットパネル型スピーカー（例えば、特許文献 3 を参照のこと）を音声出力部 221 に用いることができる。もちろん、異なるタイプのスピーカーを組み合わせたスピーカーアレイを音声出力部 2

21として用いることもできる。また、スピーカーアレイは、振動を生成する1つ以上の加振器（アクチュエータ）によって表示部219を振動させることで音声出力を行うものを含んでもよい。加振器（アクチュエータ）は、表示部219に後付けされるような形態であってもよい。図19には、ディスプレイへのパネルスピーカー技術の適用例を示している。ディスプレイ1900は、背面のスタンド1902で支持されている。ディスプレイ1900の裏面には、スピーカーユニット1901が取り付けられている。スピーカーユニット1901の左端には加振器1901-1が配置され、また、右端には加振器1901-2が配置されており、スピーカーアレイを構成している。各加振器1901-1及び1901-2が、それぞれ左右の音声信号に基づいてディスプレイ1901を振動させて音響出力することができる。スタンド1902が、低音域の音響を出力するサブウーファーを内蔵してもよい。なお、ディスプレイ1900は、有機EL素子を用いた表示部219に相当する。

[0038] 再び図2に戻って、テレビ受信装置100の構成について説明する。操作入力部222は、ユーザがテレビ受信装置100に対する操作指示の入力を行う指示入力部である。操作入力部222は、例えば、リモコン（図示しない）から送信されるコマンドを受信するリモコン受信部とボタンスイッチを並べた操作キーで構成される。また、操作入力部222は、表示部219の画面に重畳されたタッチパネルを含んでもよい。また、操作入力部222は、拡張インターフェース部205に接続されたキーボードなどの外付け入力デバイスを含んでもよい。

[0039] 拡張インターフェース部205は、テレビ受信装置100の機能を拡張するためのインターフェース群であり、例えば、アナログ映像／音声インターフェースや、USB（Universal Serial Bus）インターフェース、メモリインターフェースなどで構成される。拡張インターフェース部205は、DVI端子やHDMI（登録商標）端子やDisplay Port（登録商標）端子などからなるデジタルインターフェースを含んでい

てもよい。

[0040] 本実施形態では、拡張インターフェース205は、センサー群（後述並びに図3を参照のこと）に含まれる各種のセンサーのセンサー信号を取り込むためのインターフェースとしても利用される。センサーは、テレビ受信装置100の本体内部に装備されるセンサー、並びにテレビ受信装置100に外付け接続されるセンサーの双方を含むものとする。外付け接続されるセンサーには、テレビ受信装置100と同じ空間に存在する他のCE（Consumer Electronics）機器やIoT（Internet of Things）デバイスに内蔵されるセンサーも含まれる。拡張インターフェース205は、センサー信号をノイズ除去などの信号処理を施しさらにデジタル変換した後に取り込んでもよいし、未処理のRAWデータ（アナログ波形信号）として取り込んでもよい。

[0041] B. センシング機能

テレビ受信装置100が各種センサーを装備する目的の1つとして、ユーザが表示部219に表示された映像コンテンツを視聴する際の注視度合い（視聴質）を計測又は推定することにある。一般的には、映像コンテンツに関する満足度が高ければ、特定のシーンを注視する度合いも高くなる傾向にある。したがって、注視度を、映像コンテンツに関する「満足度」と言い換えることもできる。すなわち、本明細書中で出現する「注視度」は、同時に「満足度」と表現していることと同義であるものとする。なお、本明細書では、単に「ユーザ」という場合、特に言及しない限り、表示部219に表示された映像コンテンツを視聴する（視聴する予定がある場合も含む）視聴者のことを指すものとする。

[0042] 図3には、テレビ受信装置100に装備されるセンサー群300の構成例を示している。センサー群300は、カメラ部310と、状態センサー部320と、環境センサー部330と、機器状態センサー部340と、ユーザプロファイルセンサー部350で構成される。

[0043] カメラ部310は、表示部219に表示された映像コンテンツを視聴中の

ユーザを撮影するカメラ311と、表示部219に表示された映像コンテンツを撮影するカメラ312と、テレビ受信装置100が設置されている室内（若しくは、設置環境）を撮影するカメラ313を含む。

[0044] カメラ311は、例えば表示部219の画面の上端縁中央付近に設置され映像コンテンツを視聴中のユーザを好適に撮影する。本実施形態では、カメラ311は必須とする。

[0045] カメラ312は、例えば表示部219の画面に対向して設置され、ユーザが視聴中の映像コンテンツを撮影する。あるいは、ユーザが、カメラ312を搭載したゴーグルを装着するようにしてもよい。また、カメラ312は、映像コンテンツの音声も併せて記録（録音）する機能を備えているものとする。但し、出力する映像ストリームや音声ストリームを一時的に保持するバッファ（後述）をテレビ受信装置100内に備えている場合には、カメラ312は必須でない。

[0046] カメラ313は、例えば全天周カメラや広角カメラで構成され、テレビ受信装置100が設置されている室内（若しくは、設置環境）を撮影する。あるいは、カメラ313は、例えばロール、ピッチ、ヨーの各軸回りに回転駆動可能なカメラテーブル（雲台）に乗せたカメラであってもよい。但し、環境センサー330によって十分な環境データを取得可能な場合や環境データそのものが不要な場合には、カメラ310は不要である。

[0047] 状態センサー部320は、ユーザの状態に関する状態情報を取得する1以上のセンサーからなる。状態センサー部320は、状態情報として、例えば、ユーザの作業状態（映像コンテンツの視聴の有無）や、ユーザの行動状態（静止、歩行、走行などの移動状態、瞼の開閉状態、視線方向、瞳孔の大小）、精神状態（ユーザが映像コンテンツに没頭若しくは集中しているかなどの感動度、興奮度、覚醒度、感情や情動など）、さらには生理状態を取得することを意図している。状態センサー部320は、発汗センサー、筋電位センサー、眼電位センサー、脳波センサー、呼気センサー、ガスセンサー、イオン濃度センサー、ユーザの挙動を計測するIMU（Inertial M

Measurement Unit)、ユーザの発話を收音する音声センサー(マイクなど)などの各種のセンサーを備えていてもよい。

[0048] 環境センサー部330は、当該テレビ受信装置100が設置されている室内など環境に関する情報を計測する各種センサーからなる。例えば、温度センサー、湿度センサー、光センサー、照度センサー、気流センサー、匂いセンサー、電磁波センサー、地磁気センサー、GPS(Global Positioning System)センサー、周囲音を收音する音声センサー(マイクなど)などが環境センサー部330に含まれる。

[0049] 機器状態センサー部340は、当該テレビ受信装置100内部の状態を取得する1以上のセンサーからなる。あるいは、映像デコーダ208や音声デコーダ209などの回路コンポーネントが、入力信号の状態や入力信号の処理状況などを外部出力する機能を備えて、機器内部の状態を検出するセンサーとしての役割を果たすようにしてもよい。また、機器状態センサー部340は、当該テレビ受信装置100やその他の機器に対してユーザが行った操作を検出したり、ユーザの過去の操作履歴を保存したりするようにしてもよい。

[0050] ユーザプロフィールセンサー部350は、テレビ受信装置100で映像コンテンツを視聴するユーザに関するプロフィール情報を検出する。ユーザプロフィールセンサー部350は、必ずしもセンサー素子で構成されていなくてもよい。例えばカメラ311で撮影したユーザの顔画像や音声センサーで收音したユーザの発話などに基づいて、ユーザの年齢や性別などのユーザプロフィールを検出するようにしてもよい。また、スマートフォンなどのユーザが携帯する多機能情報端末上で取得されるユーザプロフィールを、テレビ受信装置100とスマートフォン間の連携により取得するようにしてもよい。但し、ユーザプロフィールセンサー部350は、ユーザのプライバシーや機密に関わるように機微情報まで検出する必要はない。また、同じユーザのプロフィールを、映像コンテンツの視聴の度に検出する必要はなく、一度取得したユーザプロフィール情報を例えば主制御部201内のEEPROM(

前述)に保存しておくようにしてもよい。

[0051] また、スマートフォンなどのユーザが携帯する多機能情報端末を、テレビ受信装置100とスマートフォン間の連携により、状態センサー部320若しくは環境センサー部330、ユーザプロフィールセンサー部350として活用してもよい。例えば、スマートフォンに内蔵されたセンサーで取得されるセンサー情報や、ヘルスケア機能(歩数計など)、カレンダー又はスケジュール帳・備忘録、メール、SNS(Social Network Service)といったアプリケーションで管理するデータを、ユーザの状態データや環境データに加えるようにしてもよい。

[0052] C. 注視度に基づくシーン情報推定

本実施形態に係るテレビ受信装置100は、図3に示したようなセンシング機能との組み合わせにより、ユーザの映像コンテンツに対する注視度合いを計測又は推定することができる。注視度の利用方法はさまざまである。例えば、ユーザが視聴してきたコンテンツの属性情報を集計した結果に基づいてユーザの嗜好や興味の対象を導出して、ユーザに視聴支援情報や付加価値情報を提供する情報提供装置について提案がなされている(特許文献4を参照のこと)。

[0053] ここで言うコンテンツの属性情報は、例えばコンテンツのジャンルや出演者、キーワードなどであり、例えば映像コンテンツに付随する情報である、いわゆるメタデータから属性情報を抽出することができる。また、メタデータの入手経路はさまざまであり、映像コンテンツにメタデータが重畳されるケースや、放送本編のコンテンツに付随する放送データとして配信されるケース、放送コンテンツとは異なる経路で入手されるケース(例えば、放送コンテンツのメタデータをインターネット経由で取得するケース)などさまざまである。いずれにせよ、コンテンツの属性情報は、コンテンツの制作者又は配信者側において、コンテンツ全体に対して付与されることが多い。

[0054] ユーザが興味を持ったコンテンツに関する属性情報は、ユーザへの視聴支援(例えば、自動録画予約や他のコンテンツ推薦など)、コンテンツの制作

者や配信者へのフィードバック並びにコンテンツ評価、関連商品の販売促進など、さまざまな用途が考えられる。

[0055] ところが、コンテンツ全体に関する属性情報は、一般的に、コンテンツ全体を特徴付ける情報（コンテンツのメタデータに含まれる属性情報）であり、コンテンツ中の個々のシーンを特徴付ける属性情報を示しているとは限らない。ユーザにとっては、コンテンツ全体を特徴付ける情報よりも、個々のシーンを特徴付ける情報（属性情報又はメタデータ）の方が、そのシーンに興味を持つユーザにとっては関心度の高い情報であることがある。例えば、コンテンツの制作者や配信者、広告配信者などが、コンテンツ中のシーン毎のユーザの注視度の計測結果とコンテンツ全体を特徴付ける相関関係に基づいて、視聴支援やコンテンツ評価、商品推薦などを行うと、ユーザの興味を正確に反映することができない場合がある。このために、ユーザが直接興味を持たない属性情報に従ってコンテンツや商品を推薦してしまうようなことが想定される。したがって、よりの確なコンテンツ推薦などを実現するためには、ユーザの関心を持つ個々のシーンを特徴付ける情報を取得することが必要とされる。しかしながら、コンテンツの制作者がコンテンツ制作時にすべてのシーンに異なる属性情報を添付することは困難である。

[0056] そこで、本明細書では、映像コンテンツのうちユーザが注視する特定のシーンを抽出するとともに、特定のシーンに関する情報であるメタデータ（人工知能においてニューラルネットワークを利用する場合にはラベルとも称される）を人工知能により推論して自動的に出力する技術について、以下で提案する。

[0057] 図4には、シーン取得及びシーン情報推定システム400の構成例を示している。図示のシステム400は、必要に応じて、図2に示したテレビ受信装置100内のコンポーネントや、テレビ受信装置100の外部装置（クラウド上のサーバ装置など）を用いて構成される。

[0058] 受信部401は、映像コンテンツと、映像コンテンツに付随するメタデータを受信する。映像コンテンツは、放送局（電波塔又は放送衛星など）から

送出される放送コンテンツと、OTTサービスなどのストリーム配信サーバから配信されるストリーミングコンテンツを含む。また、受信部401で受信するメタデータは、コンテンツの制作者又は配信者側においてコンテンツ全体に対して付与されるメタデータとする。そして、受信部401は、受信信号を映像ストリームと音声ストリームとメタデータに分離（デマルチプレクス）して、後段の信号処理部402とバッファ部403に出力する。

[0059] 受信部401は、例えば、テレビ受信装置100内のチューナー／復調部206、通信インターフェース部204、及びデマルチプレクサ207によって構成される。

[0060] 信号処理部402は、例えば、テレビ受信装置100内の映像デコーダ2080及び音声デコーダ209からなり、受信部401から入力した映像データストリーム及び音声データストリームをそれぞれデコードして映像情報及び音声情報を出力部404に出力する。また、信号処理部402は、デコード後の映像データストリーム及び音声データストリームをバッファ部403に出力してもよい。

[0061] 出力部404は、例えば、テレビ受信装置100内の表示部219及び音声出力部221からなり、映像情報を画面に表示出力するとともに、音声情報をスピーカーなどから音声出力する。

[0062] バッファ部403は、映像用バッファと音声用バッファを持ち、信号処理部402で復号された映像情報及び音声情報をそれぞれ一定期間だけ一時的に保持する。ここで言う一定期間とは、例えば、映像コンテンツからユーザが注視するシーンを取得するために必要な処理時間に相当する。バッファ部403は、例えば主制御部201内のRAM、あるいはその他のバッファメモリ（図示しない）であってもよい。

[0063] センサー部405は、基本的には図3に示したセンサー群300で構成される。但し、映像コンテンツを視聴中のユーザを撮影するカメラ311のみを必須とし、その他のカメラや、状態センサー320、環境センサー330などの装備は任意である。例えば、図4に示したシーン取得及びシーン情報

推定システム400はバッファ部403を備えているので、ユーザが視聴中の映像コンテンツのシーンを特定した上で記録することができる。したがって、映像コンテンツのシーンを特定するために、テレビ受信装置100の外部からテレビ受信装置100の表示部219を記録するためのカメラ312は不要である。

[0064] センサー部405は、ユーザが出力部404から出力される映像コンテンツを視聴中に、カメラ311で撮影したユーザの顔画像を注視度推定部406に出力する。また、センサー部405は、カメラ313の撮影画像や、状態センサー部320がセンシングしたユーザの状態情報、並びに環境センサー部330がセンシングした室内の環境情報なども、注視度推定部406に出力するようにしてもよい。

[0065] 注視度推定部406は、センサー部405から出力されるセンサー信号に基づいて、ユーザの視聴中の映像コンテンツに対する注視度を人工知能により推定する。本実施形態では、基本的には、注視度推定部406は、カメラ311で撮影したユーザの顔画像の認識結果に基づいて、ユーザの注視度を人工知能により推定する。例えば、注視度推定部406は、ユーザの瞳孔が開く、あるいは大きく口を開くといった顔の表情の画像認識結果に基づいて、ユーザの注視度を推定して出力する。もちろん、注視度推定部406は、カメラ311の撮影画以外のセンサー信号も入力して、ユーザの注視度を人工知能により推定するようにしてもよい。

[0066] 注視度推定部406には、人工知能の推論機能を提供するために、学習済みのニューラルネットワークを用いてもよい。図5には、注視度推定部406で利用される、注視度推定ニューラルネットワーク500の構成例を示している。注視度推定ニューラルネットワーク500は、カメラ311で撮影した画像信号やその他のセンサー信号を入力する入力層510と、中間層520と、ユーザの注視度を出力する出力層530からなる。図示の例では、中間層520は複数の中間層521、522、…からなり、ニューラルネットワーク500は深層学習(Deep Learning)を行うことがで

きる。なお、動画像や音声などの時系列情報を処理することを考慮して、中間層520において再帰結合を含むリカレントニューラルネットワーク（RNN）構造であってもよい。

[0067] 入力層510は、カメラ311で撮影した動画像ストリーム（あるいは、静止画像であってもよい）を入力ベクトルの要素に含む。基本的には、カメラ311で撮影した画像信号をRAWデータの状態のままで入力層510に入力されるものとする。

[0068] なお、カメラ311の撮影画像以外の他のセンサーのセンサー信号も注視度の計測に用いる場合には、各センサー信号に対応する入力ノードが入力層510に追加して配置される構成となる。また、画像信号や音声信号の入力などには畳み込みニューラルネットワーク（Convolutional Neural Network：CNN）を活用して特徴点の凝縮処理を行うようにしてもよい。

[0069] 出力層は、画像信号及びセンサー信号から注視度を推定した結果を、例えば0～100の連続値（若しくは、離散値）で表される注視度レベルとして出力する出力ノード530からなる。出力ノード530が連続値からなる注視度レベルを出力する場合には、出力された注視度レベルが所定値を超えたかどうかによって、ユーザが該当シーンに「注視した」又は「注視しない」と判定するようにしてもよい。あるいは、出力ノード530が「注視した」又は「注視しない」といった離散値を出力するようにしてもよい。

[0070] 注視度推定ニューラルネットワーク500の学習の過程では、顔画像やその他のセンサー信号とユーザの注視度との膨大量の組み合わせを注視度推定ニューラルネットワーク500に入力して、顔画像やその他のセンサー信号に対して尤もらしい注視度の出力ノードとの結合強度が高まるように、中間層520の各ノードの重み係数（推論係数）を更新していくことで、ユーザの顔画像（他のセンサー信号を含む場合もある）とユーザの注視度との相関関係を学習していく。そして、注視度推定ニューラルネットワーク500の利用（注視度の推定）の過程では、カメラ311で撮影された顔画像やその

他のセンサー情報を学習済みの注視度推定ニューラルネットワーク500に入力すると、ユーザの注視度が高い確度で出力される。

[0071] 図5に示すような注視度推定ニューラルネットワーク500は、例えば主制御部201内で実現される。このため、主制御部201内に、ニューラルネットワーク専用のプロセッサを含んでいてもよい。インターネット上のクラウドで注視度推定ニューラルネットワーク500を提供してもよいが、映像コンテンツに対してリアルタイムで注視度を推定していくには、注視度推定ニューラルネットワーク500はテレビ受信装置100内に配置されることが好ましい。

[0072] 例えば、エキスパート教示データベースを用いて学習を終えた注視度推定ニューラルネットワーク500を組み込んだテレビ受信装置100が出荷される。注視度推定ニューラルネットワーク500は、バックプロパゲーション（逆誤差伝播）などのアルゴリズムを利用して、継続して学習を行うようにしてもよい。あるいは、インターネット上のクラウド側で膨大なユーザから収集したデータに基づいて実施した学習結果を各家庭に設置されたテレビ受信装置100内の注視度推定ニューラルネットワーク500にアップデートすることもできるが、この点については後述に譲る。

[0073] 再び図4を参照して、シーン取得及びシーン情報推定システム400の説明を続ける。

[0074] シーン取得部407は、注視度推定部406によってユーザが注視したと判定された区間（若しくは注視度レベルが所定値を超えた区間）の映像ストリーム及び音声ストリームと、そのコンテンツ全体に関するメタデータを、バッファ部403から取得して、シーン情報推定部408に出力する。

[0075] シーン取得部407がユーザの注視度に基づいてバッファ部403から取得した映像ストリーム及び音声ストリームは、ユーザが注視していたシーンということができる。他方、コンテンツ全体に関するメタデータが、コンテンツ中の一部のシーンに対しても該当するとは限らない。何故ならば、コンテンツ全体にとってはあまり意味がないが、特定のシーンに固有の特徴（例

例えば、特定のシーンにのみ映っている物体など）に対してユーザが強い興味を持つこともあるからである。

[0076] そこで、本実施形態に係るシーン取得及びシーン情報推定システム400では、シーン情報推定部408が、ユーザが注視したシーンの映像ストリーム及び音声ストリームと、コンテンツ全体に関するメタデータとを、シーン取得部407から入力して、ユーザが注視したシーンを特徴付ける（尤もらしい）メタデータ（ラベルとも称する）が何であるかを人工知能により推定して、コンテンツ全体ではなく特定のシーンに適切な情報としてメタデータを出力するように構成されている。

[0077] シーン情報推定部408には、人工知能の推論機能を提供するために、学習済みのニューラルネットワークを用いてもよい。図6には、シーン情報推定部408で利用される、シーン情報推定ニューラルネットワーク600の構成例を示している。シーン情報推定ニューラルネットワーク600は、ユーザが注視したシーンの映像ストリーム及び音声ストリームと、コンテンツ全体のメタデータをそれぞれ入力する入力層610と、中間層620と、ユーザが注視したシーンを特徴付ける情報であるメタデータを出力する出力層630からなる。中間層620は複数の中間層621、622、…からなり、ニューラルネットワーク600は深層学習（Deep Learning）を行うことが可能であることが好ましい。なお、映像ストリームや音声ストリームなどの時系列情報を考慮して、中間層620において再帰結合を含むRNN構造であってもよい。

[0078] 入力層610は、デコードした後にバッファ部403に保持された映像ストリーム及び音声ストリームを入力ベクトルの要素に含む。但し、バッファ部403がRAWデータの状態で映像ストリーム及び音声ストリームを保持する場合には、RAWデータのまま入力層610に入力される。音声ストリームについては、時間軸方向に連続するウィンドウ毎の入力波形信号を入力ベクトルの要素に含む。前後のウィンドウ間で重複する区間があってもよい。また、ウィンドウ毎の波形信号をFFT（Fast Fourier T

r a n s f o r m) 処理した周波数信号を入力ベクトルの要素としてもよい。

[0079] また、入力層610は、コンテンツ全体を特徴付ける情報であるメタデータを入力する。コンテンツが放送番組のコンテンツの場合、メタデータは、例えば番組名や出演者名、番組内容の要約、キーワードといったテキストデータからなる。入力層610には、各テキストデータに対応する入力ノードが配置されることになる。また、画像信号や音声信号の入力などにはCNNを活用して特徴点の凝縮処理を行うようにしてもよい。

[0080] 出力層630からは、ユーザが注視したシーンを特徴付ける情報である（尤もらしい）メタデータ（ラベルとも称される）を推論できる出力が生成される。シーン毎のメタデータは、例えば番組名や出演者名、番組内容の要約、キーワードといった元のコンテンツ全体のメタデータの外、そのシーンに映り込んだ物（例えば、出演者が来ている洋服や装飾品のブランド名、出演者が手に持っているコーヒーマグの紙カップの店舗名、ロケ場所）やBGMの曲名、出演者のセリフといった元のコンテンツのメタデータにはないものが含まれることが想定される。出力層630には、これらのメタデータの各テキストデータに対応する出力ノードが配置されることになる。そして、入力層610に入力されたシーンの映像ストリーム及び音声ストリームに対して尤もらしいメタデータに該当する出力ノードが発火する。

[0081] シーン情報推定ニューラルネットワーク600の学習の過程では、ユーザが注視したシーンの映像ストリーム及び音声ストリームに対して尤もらしいメタデータの出力ノードとの結合強度が高まるように、複数層からなる中間層620の各ノードの重み係数（推論係数）を更新していくことで、シーンとメタデータとの相関関係を学習していく。そして、シーン情報推定ニューラルネットワーク600の利用、すなわちシーン情報推定の過程では、ユーザが注視したシーンの映像ストリーム及び音声ストリームに対して、尤もらしいメタデータが高い確度で出力される。

[0082] 図6に示すようなシーン情報推定ニューラルネットワーク600は、例え

ば主制御部201内で実現される。このため、主制御部201内に、ニューラルネットワーク専用のプロセッサを含んでいてもよい。例えば、エキスパート教示データベースを用いて学習を終えたシーン情報推定ニューラルネットワーク600を組み込んだテレビ受信装置100が出荷される。シーン情報推定ニューラルネットワーク600は、バックプロパゲーション（逆誤差伝播）などのアルゴリズムを利用して、学習を行うようにしてもよい。また、インターネット上のクラウド側で膨大なユーザから収集したデータに基づいて実施した学習結果を各家庭に設置されたテレビ受信装置100内のシーン情報推定ニューラルネットワーク600にアップデートすることもできるが、この点については後述に譲る。

[0083] なお、リアルタイム性が要求されない場合には、インターネット上のクラウドでシーン情報推定ニューラルネットワーク600を提供してもよい。

[0084] シーン取得部407が取得したユーザの注視シーンと、シーン情報推定部408から出力されるシーン毎に推定されるシーンの特徴付ける情報としてのメタデータ（以下、「シーンのメタデータ」とも称する）とが、シーン取得及びシーン情報推定システム400の最終出力となる。

[0085] シーン取得及びシーン情報推定システム400によるユーザの注視シーン及びそのシーンのメタデータの出力先はさまざまである。例えば、ユーザがコンテンツを視聴したテレビ受信装置100内に保存してもよいし、インターネット上のサーバにアップロードしてもよい。アップロード先のサーバは、シーン情報推定ニューラルネットワーク600が構築された人工知能サーバでもよいし、メタデータを集計するその他のサーバでもよい。また、ユーザが所持するスマートフォンなどの情報端末に、ユーザの注視シーン及びそのシーンのメタデータを出力するようにしてもよい。例えば、スマートフォン上で、シーン取得及びシーン情報推定システム400と連動するアプリケーション（仮に「コンパニオンアプリケーション」と称する）を起動させて、ユーザ本人が注視したシーンと、そのシーンのメタデータを閲覧できるようにしてもよい。

- [0086] また、ユーザが注視したシーンのメタデータの利用方法は、さまざまである。例えばユーザが視聴した映像コンテンツの評価や、他の映像コンテンツの推薦に利用することができる。また、ユーザに対してシーンに関連する商品を推薦するなどのマーケティングにも利用することができる。
- [0087] 図7には、シーン取得及びシーン情報推定システム400で実施される、映像コンテンツからシーンを取得し且つシーンを特徴付ける情報であるメタデータを出力するための処理手順をフローチャートの形式で示している。
- [0088] まず、受信部401で映像コンテンツとそのコンテンツ全体についての情報であるメタデータ（又は、メタデータに含まれる属性データ）を入力する（ステップS701）。映像コンテンツは、放送局（電波塔又は放送衛星など）から送出される放送コンテンツと、OTTサービスなどのストリーム配信サーバから配信されるストリーミングコンテンツを含む。また、受信部401で受信するメタデータは、コンテンツの制作者又は配信者側においてコンテンツ全体に対して付与されるメタデータとする。
- [0089] 次に、受信部401で入力したコンテンツの映像ストリーム及び音声ストリームを信号処理部402で処理して出力部404で出力するとともに、コンテンツ及びメタデータのバッファリングを行う（ステップS702）。
- [0090] そして、出力部404で提示するコンテンツを視聴中におけるユーザの注視度を、注視度推定部406で人工知能により推定する（ステップS703）。
- [0091] 注視度推定部406は、基本的には、カメラ311で撮影したユーザの顔画像の認識結果に基づいて、ユーザの注視度を計測し又は人工知能により推定する。注視度推定部406は、カメラ311の撮影画像以外の他のセンサーのセンサー信号も注視度の計測に用いることもできる。また、注視度推定部406には、人工知能による推定を行うために、顔画像やその他のセンサー情報とユーザの注視度との相関関係を学習した注視度推定ニューラルネットワーク500（図5を参照のこと）が用いられる。
- [0092] 次に、シーン取得部407は、注視度推定部406によってユーザが注

視したと判定されると（ステップS70のYes）、ユーザが注視した区間の映像ストリーム及び音声ストリームを、ユーザが注視したシーンとしてバッファ部403から取得するとともに、そのコンテンツ全体のメタデータをバッファ部403から取得して、シーン情報推定部408に出力する（ステップS705）。

[0093] 次いで、シーン情報推定部408は、ユーザが注視したシーンの映像ストリーム及び音声ストリームと、コンテンツ全体に関するメタデータとを、シーン取得部407から入力して、ユーザが注視したシーンを特徴付ける（尤もらしい）メタデータが何であるかを人工知能により推定する（ステップS706）。シーン情報推定部408には、人工知能による推定を行うために、ユーザが注視したシーンと、シーンに対して尤もらしいメタデータとの相関関係を学習したシーン情報推定ニューラルネットワーク600（図6を参照のこと）が用いられる。

[0094] そして、ユーザが注視したシーンと、そのシーンを特徴付ける（尤もらしい）、シーンのメタデータとを所定の出力先に出力して（ステップS707）、本処理を終了する。

[0095] 図8には、シーン取得及びシーン情報推定システム800の変形例を示している。図示のシステム800は、必要に応じて、図2に示したテレビ受信装置100内のコンポーネントや、テレビ受信装置100の外部装置（クラウド上のサーバ装置など）を用いて構成される。

[0096] 受信部801、信号処理部802、バッファ部803、出力部804、及びセンサー部805の構成及び動作は、図4に示したシーン情報推定システム400と同様なので、ここでは詳細な説明を省略する。

[0097] シーン情報推定部806は、図4に示したシーン情報推定システム400における注視度推定部406、シーン取得部407、及びシーン情報推定部408が行う処理を、単一のニューラルネットワークで一括して実現する。

[0098] 図9には、シーン情報推定部806で利用されるシーン情報推定ニューラルネットワーク900の構成例を示している。シーン情報推定ニューラルネ

ットワーク900は、入力層910と、中間層920と、出力層930からなる。図示の例では、中間層920は、複数の中間層921、922、…からなり、シーン情報推定ニューラルネットワーク900は深層学習（DL）を行うことができる。なお、動画像や音声などの時系列情報を考慮して、中間層920において再帰結合を含むリカレントニューラルネットワーク（RNN）構造であってもよい。

[0099] 入力層910は、カメラ311で撮影した画像信号やその他のセンサー情報、バッファ部803に一時的に保持される映像ストリーム及び音声ストリーム、並びにコンテンツのメタデータを入力するノードを備えている。

[0100] 入力層910は、カメラ311で撮影した動画像ストリーム（あるいは、静止画像であってもよい）や、コンテンツの映像ストリーム及び音声ストリームなどを入力ベクトルの要素に含む。基本的には、画像信号や音声信号をRAWデータの状態のまま入力層910に入力されるものとする。音声ストリームについては、時間軸方向に連続するウィンドウ毎の入力波形信号を、入力層910の各ノードに入力する入力ベクトルとする。前後のウィンドウ間で重複する区間があってもよい。また、ウィンドウ毎の波形信号をFFT処理した周波数信号を入力ベクトルとしてもよい（同上）。また、画像信号や音声信号の入力などには畳み込みニューラルネットワークCNNを活用して特徴点の凝縮処理を行うようにしてもよい。

[0101] また、放送コンテンツの番組名や出演者名、番組内容の要約、キーワードといったテキストデータからなるメタデータを入力層910に入力する場合、入力層910には、各テキストデータに対応する入力ノードが配置されることになる。

[0102] 一方、出力層930からは、ユーザの注視したシーンと、そのシーンの特徴付ける（尤もらしい）メタデータ（ラベルとも称される）を推定できる出力が生成される。シーン毎のメタデータは、例えば番組名や出演者名、番組内容の要約、キーワードといった元のコンテンツ全体のメタデータその他、そのシーンに映り込んだ物（例えば、出演者が来ている洋服や装飾品のブラン

ド名、出演者が手に持っているコーヒーの紙カップの店舗名、ロケ場所) や BGMの曲名、出演者のセリフといった元のコンテンツのメタデータにはないものが含まれることが想定される。したがって、出力層930には、これらのメタデータの各テキストデータに対応する出力ノードが配置されることになる。そして、入力層910に入力された映像ストリーム及び音声ストリームに対して、ユーザの注視度が高いシーンの映像ストリーム及び音声ストリームと、そのシーンに対して尤もらしいメタデータに該当する出力ノードが発火する。

[0103] シーン情報推定ニューラルネットワーク900の学習の過程では、顔画像やその他のセンサー信号と、時々刻々と入力される映像ストリーム及び音声ストリームに対して、尤もらしい注視シーン並びにそのシーンに対して尤もらしいメタデータの出力ノードとの結合強度が高まるように、複数層からなる中間層920の各ノードの重み係数(推論係数)を更新していくことで、シーンとメタデータとの相関関係を学習していく。そして、シーン情報推定ニューラルネットワーク900の利用、すなわちシーン情報推定の過程では、顔画像やその他のセンサー信号と、時々刻々と入力される映像ストリーム及び音声ストリームに対して、尤もらしい注視シーン並びにそのシーンに対して尤もらしいメタデータが高い確度で出力される。

[0104] 図9に示すようなシーン情報推定ニューラルネットワーク900は、例えば主制御部201内で実現される。このため、主制御部201内又は主制御部201とは別の処理回路内に、ニューラルネットワーク専用のプロセッサを含んでいてもよい。例えば、エキスパート教示データベースを用いて学習を終えたシーン情報推定ニューラルネットワーク900を組み込んだテレビ受信装置100が出荷される。シーン情報推定ニューラルネットワーク900は、バックプロパゲーション(逆誤差伝播)などのアルゴリズムを利用して、学習を行うようにしてもよい。また、インターネット上のクラウド側で膨大なユーザから収集したデータに基づいて実施した学習結果を各家庭に設置されたテレビ受信装置100内のシーン情報推定ニューラルネットワーク

900にアップデートすることもできるが、この点については後述に譲る。

[0105] なお、リアルタイム性が要求されない場合には、インターネット上のクラウドでシーン情報推定ニューラルネットワーク900を提供してもよい。

[0106] シーン情報推定部806によって特定される、尤もらしいユーザの注視シーンと、そのシーンに対して尤もらしいメタデータが、シーン取得及びシーン情報推定システム800の最終出力となる。

[0107] D. 注視度に基づくシーン情報推定結果のフィードバック

ここでは、シーン取得及びシーン情報推定システム400によるシーン情報推定結果のユーザへのフィードバック方法について、シーン情報推定システム400の動作とともに説明する。

[0108] 図10に示すように、テレビ受信装置100の画面に映像（放送コンテンツやOTTサービスのストリーミング動画など）を表示中に、例えば画面の上端縁の中央付近に設置されたカメラ311によって、ユーザを撮影し続ける。

[0109] 注視度推定部406は、カメラ311の撮影画像から顔認識して、ユーザの注視度を計測する。図11には、カメラ311の撮影画像から2人のユーザの顔画像1101及び1102を認識している様子を示している。注視度推定部406は、顔画像1101及び1102のうち少なくとも一方について注視度推定を実施する。ユーザのコンテンツに対する注視度の人工知能による推定には、図5に示した注視度推定ニューラルネットワーク500を用いてもよい。注視度推定ニューラルネットワーク500は、カメラ311の撮影画像以外に、他のセンサーのセンサー信号に基づいて、ユーザの注視度を推定するようにしてもよい。

[0110] カメラ311の撮影画像から認識された顔画像1101又は1102が注視していることを示している、と注視度推定部406が推定した場合には、注視していると推定される区間においてバッファ部403に保持されている映像ストリーム及び音声ストリームを、ユーザの注視シーンとして捕捉する。図12には、ユーザの注視度の計測に伴って注視シーンを捕捉する様子を

例示している。

[0111] シーン取得及びシーン情報推定システム400が出力するユーザの注視シーン及びそのシーンを特徴付ける情報であるメタデータは、インターネット上のサーバにアップロードされる。例えば、コンテンツ推薦サーバは、ユーザの注視シーン及びそのシーンのメタデータに基づいて、そのユーザに視聴を薦める類似又は関連コンテンツを検索して、ユーザに推薦コンテンツの情報を提供する。この種のコンテンツ推薦サーバは、CF (Collaborative Filtering) やCBF (Content Based Filtering) といったアルゴリズムを用いて、ユーザの注視シーンから類似又は関連コンテンツを検索するようにしてもよい。あるいは、コンテンツ推薦サーバは、前述した人工知能サーバとして、メタデータとコンテンツとの相関関係を学習したニューラルネットワークを用いて類似又は関連コンテンツを抽出するようにしてもよい。

[0112] コンテンツ推薦サーバは、いずれのアルゴリズムに従って類似又は関連コンテンツを検索するにせよ、コンテンツ全体ではなく、コンテンツ中でユーザが注視した特定のシーンのメタデータに基づいて類似又は関連コンテンツの検索を行うことになる。このため、番組名や出演者といったコンテンツ全体に関わるメタデータとは関連が低いが、ユーザが注視したシーンにしか映っていないような細かな事象（例えば、出演者が来ている服やサングラス・時計といった装飾品のブランド）に関連するようなコンテンツをユーザに推薦することができる。

[0113] 図13には、ユーザが過去に注視したシーンをフィードバックする画面の構成例を示している。図示の画面上には、テレビ受信装置100で映像コンテンツを視聴してきたユーザがこれまでに注視したことのあるシーンの代表画像の一覧がマトリックス状に表示されている。ここで言う「代表画像」は、例えばコンテンツから切り出したシーンの先頭の画像でもよいし、シーンの中から所定のアルゴリズム又はランダムにキャプチャーした画像でもよい。画面上の表示形式は、マトリックス状に限られるものではなく、横方向に

一列、又は縦方向に一列で、仮想的に無限長のリングを構成するようにしてもよい。また、3次元空間の任意に位置に配置された多くのシーンを、当該3次元空間を所定の視点から視た場合に視覚的に表示される2次元平面として表示し、視点を変更できるように構成されてもよい。

[0114] ユーザは、例えばリモコンの十字キーなどを使って、一覧のうち気になるシーンの代表画像を選択することができる。そして、いずれかの代表画像を選択中に、例えばリモコンの決定ボタンを押すと、そのシーンの選択が確定する。なお、代表画像の一覧から特定のシーンを選択するための上記のリモコン操作は、例えば人工知能の機能を持つAI (Artificial Intelligence) スピーカーなどの音声エージェント機能を通じてマイクから音声指示することもできる。あるいは、ジェスチャーによってリモコン操作を行ってもよい。ジェスチャーは、手などに限らず頭部を含めた身体のすべての部位を使うことができる。また、スマートグラスなどを装着している場合には、視線を認識して、ウィンクなどの目を用いたジェスチャーにより選択確定指示を行うようにしてもよい。さらに、モーションセンサを装備して、人工知能の機能により頭部の動きを感知できるワイヤレスヘッドホンを用いてもよい。ワイヤレスヘッドホンなどの無線機器には、タッチセンサを用いて選択確定指示を行うようにしてもよい。

[0115] 図14には、ある1つの代表画像を選択した後に遷移する画面の構成例を示している。図示の例では、代表画像の一覧(図13を参照のこと)の中からユーザが選択した1つの代表画像が拡大表示されるとともに、そのシーンに対してシーン情報推定システム400によって推定された情報であるシーンのメタデータの一部又は全部(MD#1、MD#2、…)が表示される。代表画像を表示した際、静止画ではなく、ユーザが注視した区間の映像ストリーム及び音声ストリームを再生出力するようにしてもよい。

[0116] ユーザは、注視シーンの代表画像又は動画像を視聴することで、過去に注視したシーンの記憶を思い起こすことができる。また、ユーザは、自分が注視したシーンの代表画像と一緒にメタデータを見ることで、自分がそのシー

ンに注視した理由若しくは根拠を確認することができる。そして、ユーザは、シーン情報推定部408が人工知能により自動的にそのシーンに対して付与した、そのシーンを特徴付けるメタデータを変更したい場合などには、リモコン操作や音声エージェント機能などを通じて、メタデータの編集（変更、追加、削除など）を行い、個人用にシーンに関連付けられたメタデータをカスタマイズすることができる。

[0117] さらにユーザは、リモコン操作や音声エージェントを通じて、選択した注視シーンに基づいて推薦される類似又は関連コンテンツの提示を指示することができる。図15には、選択したシーンのメタデータに基づいて推薦される類似又は関連コンテンツを提示する画面の構成例を示している。図示の画面の左半分には、ユーザの注視シーンがそのシーンのメタデータとともに表示されている。また、画面の右半分には、類似又は関連コンテンツのリストが表示される。類似又は関連コンテンツのリストには、各々の類似又は関連コンテンツの代表画像とメタデータが表示される。ユーザは、代表画像やメタデータに基づいて、各々の類似又は関連コンテンツを見たいかどうかを判断することができる。そして、ユーザは、視聴を希望する類似又は関連コンテンツが見つかったときには、リモコン操作や音声エージェント機能などを通じて、そのコンテンツの再生を指示して、視聴を開始することができる。

[0118] また、ユーザは、テレビ受信装置100の大画面だけではなく、スマートフォンなどの情報端末の小さな画面上でも、注視シーンのフィードバックを表示することができる。例えば、スマートフォン上で、シーン取得及びシーン情報推定システム400と連動するコンパニオンアプリケーションを起動させて、ユーザ本人が注視したシーンと、そのシーンのメタデータを閲覧できるようにしてもよい。

[0119] 図16には、ユーザが注視したシーンをスマートフォンの画面上でフィードバックする様子を例示している。スマートフォンの画面は小さく、図13に示したテレビ画面のように複数のシーンの代表画像を一覧表示することはできない。このため、一画面には1つのシーンの代表画像のみを表示する。

但し、複数のシーンの代表画像を、例えば図17に示すように仮想的にカルーセル状に配置したり、あるいは図18に示すように仮想的にマトリックス状に配置したりしておき、ユーザがスマートフォンの画面に表示された代表画像を水平方向又は垂直方向にフリック操作したことに応じて、カルーセル又はマトリックス上でフリック操作した方向に隣接する代表画像に遷移させるようにしてもよい。

[0120] 再び図16に戻って、スマートフォン上のフィードバック画面の構成について説明する。代表画像を画面に表示した際、静止画ではなく、ユーザが注視した区間の映像ストリーム及び音声ストリームを再生出力するようにしてもよい。また、代表画像の下に、そのシーンに対してシーン情報推定システム400によって推定されたメタデータの一部又は全部が表示される。

[0121] 本明細書で開示するシーン情報推定システム400によれば、ユーザは、注視シーンの代表画像又は動画像を視聴することで、過去に注視したシーンの記憶を思い起こすことができる、という効果を有する。また、本明細書で開示するシーン情報推定システム400によれば、ユーザは、自分が注視したシーンの代表画像と一緒にメタデータを見ることで、自分がそのシーンに注視した理由若しくは根拠を確認することができる、という効果を有する。さらに、本明細書で開示するシーン情報推定システム400によれば、ユーザは、そのシーンに対して付与されたメタデータを変更したい場合などには、タッチパネルを利用したスマートフォンの編集機能を利用して、メタデータの編集（変更、追加、削除など）を行うことができる、という効果を有する。またさらに、本明細書で開示するシーン情報推定システム400によれば、ユーザの注視度の高いシーンを集積して学習することにより、ユーザの注視度すなわち「満足度」の高い映像コンテンツの傾向を学習することができる。したがって、人工知能により、学習機能を通じて、「満足度」の高い映像コンテンツをユーザに提供できる可能性を高めることができる、という効果を有する。

[0122] E. ニューラルネットワークのアップデート

これまで、映像コンテンツ中でユーザが注視したシーンに関するメタデータを人工知能により推定する過程で用いられる、注視度推定用ニューラルネットワーク500、シーン情報推定ニューラルネットワーク600、並びにシーン情報推定ニューラルネットワーク900について説明してきた。

[0123] これらのニューラルネットワークは人工知能の機能として用いられ、各家庭に設置されたテレビ受信装置100というユーザが直接操作することができる間装置又はその装置が設置された例えば家庭のような動作環境（以下、「ローカル環境」とも呼ぶ）で動作する。人工知能の機能としてニューラルネットワークをローカル環境で動作させることの効果の1つは、例えば、これらのニューラルネットワークに対してバックプロパゲーション（逆誤差伝播）などのアルゴリズムを利用し、ユーザからのフィードバックなどを教師データとして学習を行うことを容易にリアルタイムに実現できることである。ユーザからのフィードバックは、例えば、注視度推定ニューラルネットワーク500が推定したユーザの注視度や、シーン情報推定ニューラルネットワーク600又は900が推定したシーンを特徴付ける情報であるメタデータに対するユーザの評価である。ユーザフィードバックは、例えばOK（良）、NG（不良）といった簡単なものでもよい。ユーザフィードバックは、例えば操作入力部222やリモコン、人工知能の一形態である音声エージェント、連携するスマートフォンなどを介してテレビ受信装置100に入力される。したがって、これらの人工知能の機能としてニューラルネットワークをローカル環境で動作させることの効果の別の側面は、ユーザフィードバックを利用した学習により、ニューラルネットワークを特定のユーザにカスタマイズ若しくはパーソナライズすることができることである。

[0124] 他方、インターネット上のサーバ装置の集合体であるクラウド上で動作する1つ以上のサーバ装置（以下、単に「クラウド」とも呼ぶ）において、膨大な数のユーザからデータを収集して、人工知能の機能としてニューラルネットワークの学習を積み重ね、その学習結果を用いて各家庭のテレビ受信装置100内のニューラルネットワークをアップデートする方法も考えられる。

クラウドで人工知能の機能を果たすニューラルネットワークのアップデートを行うことの効果の1つは、大量のデータで学習することにより、より確度の高いニューラルネットワークを構築することができることである。

[0125] 図20には、クラウドを利用した人工知能システム2000の構成例を模式的に示している。図示のクラウドを利用した人工知能システム2000は、ローカル環境2010とクラウド2020からなる。

[0126] ローカル環境2010は、テレビ受信装置100を設置した動作環境（家庭）、若しくは家庭内に設置されたテレビ受信装置100に相当する。図20には、簡素化のため1つのローカル環境2010しか描いていないが、実際には、1つのクラウド2020に対して膨大な数のローカル環境が接続されることが想定される。また、本実施例では、ローカル環境2010としてテレビ受信装置100又はテレビ受信装置100の動作する家庭のような動作環境を主に例示したが、ローカル環境2010は、スマートフォンやウェアラブルデバイスなど、ユーザが直接操作することができる任意の装置又は装置の動作する環境（駅、バス停、空港、ショッピングセンターのような公共施設、工場や職場などの労働設備を含む）であればよい。

[0127] 上述したように、テレビ受信装置100内には、人工知能として、注視度推定ニューラルネットワーク500とシーン情報推定ニューラルネットワーク600、又はシーン情報推定ニューラルネットワーク900が配置されている。テレビ受信装置100内に搭載され、実際に利用に供されるこれらのニューラルネットワークのことを、ここでは運用ニューラルネットワーク2011と総称することにする。運用ニューラルネットワーク2011は、膨大なサンプルデータからなるエキスパート教示データベースを用いて、事前に学習が行われていることを想定している。

[0128] 一方、クラウド2020には、人工知能機能を提供する人工知能サーバ（前述）（1つ以上のサーバ装置から構成される）が装備されている。人工知能サーバは、運用ニューラルネットワーク2021と、運用ニューラルネットワーク2022を評価する評価ニューラルネットワーク2022が配設さ

れている。運用ニューラルネットワーク2021は、ローカル環境2010に配置された運用ニューラルネットワーク2011と同一構成であり、膨大なサンプルデータからなるエキスパート教示データベースを用いて、事前に学習が行われていることを想定している。また、評価ニューラルネットワーク2022は、運用ニューラルネットワーク2021の学習状況の評価に用いられるニューラルネットワークである。

- [0129] ローカル環境2010側では、運用ニューラルネットワーク2011は、カメラ311の撮影画像などのセンサー情報とユーザプロフィールを入力して、ユーザプロフィールに適合した注視度やシーン毎のメタデータを出力する。但し、運用ニューラルネットワーク2011がシーン情報推定ニューラルネットワーク600の場合、ユーザが注視したシーンの映像ストリームと元のコンテンツのメタデータを入力とする。ここでは、簡素化のため、運用ニューラルネットワーク2011への入力を単に「入力値」と呼び、運用ニューラルネットワーク2012からの出力を単に「出力値」と呼ぶことにする。
- [0130] ローカル環境2010のユーザ（例えば、テレビ受信装置100の視聴者）は、運用ニューラルネットワーク2011の出力値を評価して、例えば操作入力部222やリモコン、音声エージェント、連携するスマートフォンなどを介してテレビ受信装置100に評価結果をフィードバックする。ここでは、説明の簡素化のため、ユーザフィードバックは、OK（0）又はNG（1）のいずれかであるとする。
- [0131] ローカル環境2010からクラウド2020へ、運用ニューラルネットワーク2011の入力値と出力値、及びユーザフィードバックの組み合わせからなるフィードバックデータがクラウド2020に送信される。クラウド2020内では、膨大な数のローカル環境から送られてきたフィードバックデータが、フィードバックデータベース2023に蓄積されていく。フィードバックデータベース2023には、運用ニューラルネットワーク2011の入力値及び出力値とユーザとの対応関係を記述した膨大量のフィードバックデ

ータが蓄積される。

- [0132] また、クラウド2020は、運用ニューラルネットワーク2011の事前学習に用いられた、膨大なサンプルデータからなるエキスパート教示データベース2024を所有し又は利用が可能である。個々のサンプルデータは、センサー情報及びユーザプロファイルと運用ニューラルネットワーク2011（若しくは、2021）の出力値との対応関係を記述した教師データである。
- [0133] フィードバックデータベース2023からフィードバックデータを取り出すと、フィードバックデータに含まれる入力値（例えば、センサー情報とユーザプロファイルの組み合わせ）が運用ニューラルネットワーク2021に入力される。また、評価ニューラルネットワーク2022には、運用ニューラルネットワーク2021の出力値と、対応するフィードバックデータに含まれる入力値（例えば、センサー情報とユーザプロファイルの組み合わせ）が入力され、評価ニューラルネットワーク2022はユーザフィードバックを出力する。
- [0134] クラウド2020内では、第1ステップとしての評価ニューラルネットワーク2022の学習と、第2ステップとしての運用ニューラルネットワーク2021の学習が交互に実施される。
- [0135] 評価ニューラルネットワーク2022は、運用ニューラルネットワーク2021への入力値と、運用ニューラルネットワーク2021の出力に対するユーザフィードバックとの対応関係を学習するネットワークである。したがって、第1ステップでは、評価ニューラルネットワーク2022は、運用ニューラルネットワーク2021の出力値と、対応するフィードバックデータに含まれるユーザフィードバックとを入力して、運用ニューラルネットワーク2021の出力値に対して自身が出力するユーザフィードバックが、運用ニューラルネットワーク2021の出力値に対する現実のユーザフィードバックと一致するように学習する。この結果、評価ニューラルネットワーク2022は、運用ニューラルネットワーク2021の出力に対して、現実のユ

ーザと同じようなユーザフィードバック（OK又はNG）を出力するように、学習されていく。

[0136] 続く第2ステップでは、評価ニューラルネットワーク2022を固定して、今度は運用ニューラルネットワーク2021の学習を実施する。上述したように、フィードバックデータベース2023からフィードバックデータを取り出すと、フィードバックデータに含まれる入力値が運用ニューラルネットワーク2021に入力され、評価ニューラルネットワーク2022には、運用ニューラルネットワーク2021の出力値と、対応するフィードバックデータに含まれるユーザフィードバックのデータが入力され、評価ニューラルネットワーク2022は現実のユーザと等しいユーザフィードバックを出力する。

[0137] このとき、運用ニューラルネットワーク2021は、ニューラルネットワークの出力層からの出力に対して評価関数（例えば、ロス関数）を適用して、その値が最小となるようにバックプロパゲーションを用いて学習を実施する。例えば、ユーザフィードバックを教師データとする場合、運用ニューラルネットワーク2021は、すべての入力値に対して評価ニューラルネットワーク2022の出力がOK（0）となるように学習する。このような学習を実施することによって、運用ニューラルネットワーク2021は、いかなる入力値（センサー情報、ユーザプロフィールなど）に対しても、ユーザがOKとフィードバックする出力値（注視度、シーンに対するメタデータなど）を出力することができるようになる。

[0138] また、運用ニューラルネットワーク2021の学習時において、エキスパート教示データベース2024を教師データに用いてもよい。また、ユーザフィードバックやエキスパート教示データベース2024など、2以上の教師データを用いて学習を行うようにしてもよい。この場合、教師データ毎に算出したロス関数を重み付け加算して、最小となるように運用ニューラルネットワーク2021の学習を行うようにしてもよい。

[0139] 上述したような第1ステップとしての評価ニューラルネットワーク202

2の学習と、第2ステップとしての運用ニューラルネットワーク2021の学習が交互に実施することによって、運用ニューラルネットワーク2021の確度が向上していく。そして、学習により確度が向上した運用ニューラルネットワーク2021における推論係数を、ローカル環境2010における運用ニューラルネットワーク2011に提供することで、ユーザもさらに学習が進んだ運用ニューラルネットワーク2011を享受することができる。

[0140] 例えば、運用ニューラルネットワーク2011の推論係数のビットストリームを圧縮して、クラウド2020からローカル環境へダウンロードすればよい。圧縮してもビットストリームのサイズが大きいときには、層毎あるいは領域毎に推論係数を分割して、複数回に分けて圧縮ビットストリームをダウンロードするようにしてもよい。

産業上の利用可能性

[0141] 以上、特定の実施形態を参照しながら、本明細書で開示する技術について詳細に説明してきた。しかしながら、本明細書で開示する技術の要旨を逸脱しない範囲で当業者が該実施形態の修正や代用を成し得ることは自明である。

[0142] 本明細書では、本明細書で開示する技術をテレビ受信機に適用した実施形態を中心に説明してきたが、本明細書で開示する技術の要旨はこれに限定されるものではない。映像や音声などさまざまな再生コンテンツをユーザに提示するさまざまなタイプのコンテンツ再生装置にも、同様に本明細書で開示する技術を適用することができる。

[0143] 要するに、例示という形態により本明細書で開示する技術について説明してきたのであり、本明細書の記載内容を限定的に解釈するべきではない。本明細書で開示する技術の要旨を判断するためには、特許請求の範囲を参酌すべきである。

[0144] なお、本明細書の開示の技術は、以下のような構成をとることも可能である。

[0145] (1) コンテンツを視聴しているユーザの注視度をセンサー情報に基づいて

人工知能により推定する注視度推定部と、

前記注視度推定部の推定結果に基づいて、前記コンテンツ中で前記ユーザが注視するシーンの映像と前記コンテンツに関する情報を取得する取得部と

、

前記ユーザが注視するシーンの映像と前記コンテンツに関する情報に基づいて、前記ユーザが注視しているシーンに関する情報を人工知能により推定するシーン情報推定部と、

を具備する人工知能情報処理装置。

[0146] (2) 前記シーン情報推定部は、人工知能による推定として、シーンの映像と前記コンテンツに関する情報と、シーンに関する情報との相関関係を学習したニューラルネットワークを利用して、前記ユーザが注視するシーンと相関関係のある情報を推定する、
上記(1)に記載の人工知能情報処理装置。

[0147] (3) 前記注視度推定部は、人工知能による推定として、センサー情報と前記ユーザの注視度との相関関係を学習したニューラルネットワークを利用して、前記コンテンツを視聴している前記ユーザに関するセンサー情報と相関関係のある注視度合いを推定する、
上記(1)又は(2)のいずれかに記載の人工知能情報処理装置。

[0148] (4) 前記センサー情報は、前記コンテンツを視聴している前記ユーザをカメラで撮影した撮影画像を少なくとも含み、
前記注視度推定部は、人工知能による推定として、顔認識結果と前記ユーザの注視度との相関関係を学習したニューラルネットワークを利用して、前記撮影画像に顔認識結果と相関関係のある注視度を推定する、
上記(3)に記載の人工知能情報処理装置。

[0149] (5) 前記コンテンツは、放送コンテンツ又はストリーム配信コンテンツの少なくとも一方を含む、
上記(1)乃至(4)のいずれかに記載の人工知能情報処理装置。

[0150] (6) 前記ユーザが注視しているシーンに関する情報を外部に出力する、

上記（１）乃至（５）のいずれかに記載の人工知能情報処理装置。

[0151] （７）前記ユーザが注視しているシーン及びそのシーンに関する情報をユーザに提示する、

上記（１）乃至（６）のいずれかに記載の人工知能情報処理装置。

[0152] （８）コンテンツを視聴しているユーザに関するセンサー情報を入力する入力部と、

センサー情報、コンテンツ及びコンテンツの情報と、ユーザが注視するシーンに関する情報との相関関係を学習したニューラルネットワークを利用して、前記ユーザが注視するシーンと相関関係のある情報を推定するシーン情報推定部と、

を具備する人工知能情報処理装置。

[0153] （９）前記センサー情報は、前記コンテンツを視聴している前記ユーザをカメラで撮影した撮影画像を少なくとも含み、

前記シーン情報推定部は、顔画像、コンテンツ及びコンテンツの情報と、ユーザが注視するシーンに関する情報との相関関係を学習したニューラルネットワークを利用して、前記ユーザが注視するシーンと相関関係のある情報を推定する、

上記（８）に記載の人工知能情報処理装置。

[0154] （１０）コンテンツを視聴しているユーザの注視度をセンサー情報に基づいて人工知能により推定する注視度推定ステップと、

前記注視度推定ステップにおける推定結果に基づいて、前記コンテンツ中で前記ユーザが注視するシーンの映像と前記コンテンツに関する情報を取得する取得ステップと、

前記ユーザが注視するシーンの映像と前記コンテンツに関する情報に基づいて、前記ユーザが注視しているシーンに関する情報を人工知能により推定するシーン情報推定ステップと、

を有する人工知能情報処理方法。

符号の説明

- [0155] 100…テレビ受信装置、201…制御部、202…バス
203…ストレージ部、204…通信インターフェース（I F）部
205…拡張インターフェース（I F）部
206…チューナー／復調部、207…デマルチプレクサ
208…映像デコーダ、209…音声デコーダ
210…文字スーパーデコーダ、211…字幕デコーダ
212…字幕合成部、213…データデコーダ、214…キャッシュ部
215…アプリケーション（A P）制御部、216…ブラウザ部
217…音源部、218…映像合成部、219…表示部
220…音声合成部、221…音声出力部、222…操作入力部
300…センサー群、310…カメラ部、311～313…カメラ
320…状態センサー部、330…環境センサー部
340…機器状態センサー部、350…ユーザプロファイルセンサー部
400…シーン取得及びシーン情報推定システム、401…受信部
402…信号処理部、403…バッファ部、404…出力部
405…センサー部、406…注視度推定部、407…シーン取得部
408…シーン情報推定部
500…注視度推定ニューラルネットワーク
510…入力層、520…中間層、530…出力層
600…シーン情報推定ニューラルネットワーク
610…入力層、620…中間層、630…出力層
800…シーン取得及びシーン情報推定システム、801…受信部
802…信号処理部、803…バッファ部、804…出力部
805…センサー部、806…シーン情報推定部
900…シーン情報推定ニューラルネットワーク
910…入力層、920…中間層、930…出力層
1900…ディスプレイ、1901…スピーカユニット
1902…スタンド、1901-1、1902-2…加振器

2000…クラウドを利用した人工知能システム

2010…ローカル環境、2011…運用ニューラルネットワーク

2020…クラウド、2021…運用ニューラルネットワーク

2022…評価ニューラルネットワーク

2023…フィードバックデータベース

2024…エキスパート教示データベース

請求の範囲

- [請求項1] コンテンツを視聴しているユーザの注視度をセンサー情報に基づいて人工知能により推定する注視度推定部と、
- 前記注視度推定部の推定結果に基づいて、前記コンテンツ中で前記ユーザが注視するシーンの映像と前記コンテンツに関する情報を取得する取得部と、
- 前記ユーザが注視するシーンの映像と前記コンテンツに関する情報に基づいて、前記ユーザが注視しているシーンに関する情報を人工知能により推定するシーン情報推定部と、
- を具備する人工知能情報処理装置。
- [請求項2] 前記シーン情報推定部は、人工知能による推定として、シーンの映像と前記コンテンツに関する情報と、シーンに関する情報との相関関係を学習したニューラルネットワークを利用して、前記ユーザが注視するシーンと相関関係のある情報を推定する、
- 請求項1に記載の人工知能情報処理装置。
- [請求項3] 前記注視度推定部は、人工知能による推定として、センサー情報と前記ユーザの注視度との相関関係を学習したニューラルネットワークを利用して、前記コンテンツを視聴している前記ユーザに関するセンサー情報と相関関係のある注視度合いを推定する、
- 請求項1に記載の人工知能情報処理装置。
- [請求項4] 前記センサー情報は、前記コンテンツを視聴している前記ユーザをカメラで撮影した撮影画像を少なくとも含み、
- 前記注視度推定部は、人工知能による推定として、顔認識結果と前記ユーザの注視度との相関関係を学習したニューラルネットワークを利用して、前記撮影画像に顔認識結果と相関関係のある注視度を推定する、
- 請求項3に記載の人工知能情報処理装置。
- [請求項5] 前記コンテンツは、放送コンテンツ又はストリーム配信コンテンツ

の少なくとも一方を含む、

請求項 1 に記載の人工知能情報処理装置。

[請求項6] 前記ユーザが注視しているシーンに関する情報を外部に出力する、
請求項 1 に記載の人工知能情報処理装置。

[請求項7] 前記ユーザが注視しているシーン及びそのシーンに関する情報をユーザに提示する、
請求項 1 に記載の人工知能情報処理装置。

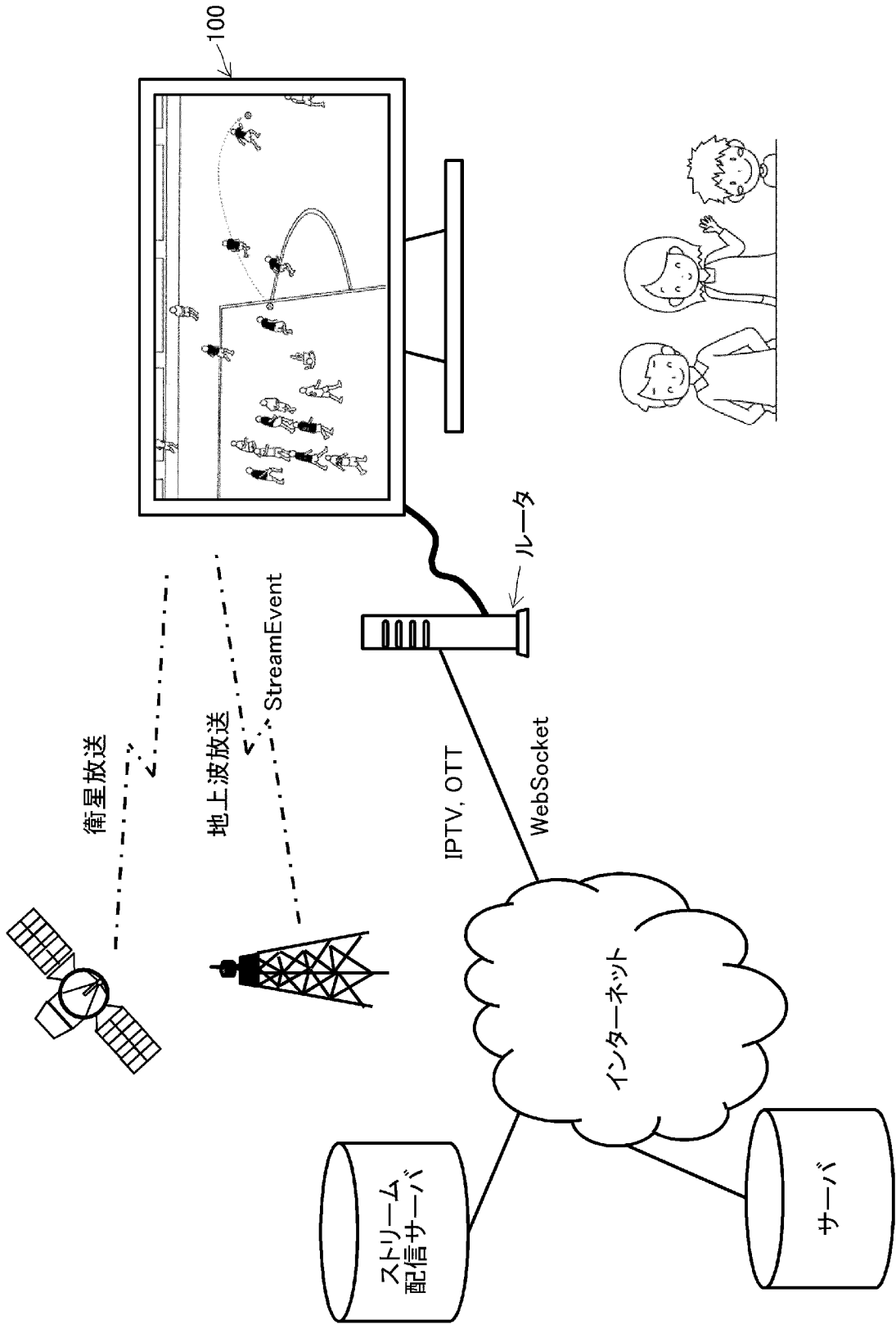
[請求項8] コンテンツを視聴しているユーザに関するセンサー情報を入力する入力部と、
センサー情報、コンテンツ及びコンテンツの情報と、ユーザが注視するシーンに関する情報との相関関係を学習したニューラルネットワークを利用して、前記ユーザが注視するシーンと相関関係のある情報を推定するシーン情報推定部と、
を具備する人工知能情報処理装置。

[請求項9] 前記センサー情報は、前記コンテンツを視聴している前記ユーザをカメラで撮影した撮影画像を少なくとも含み、
前記シーン情報推定部は、顔画像、コンテンツ及びコンテンツの情報と、ユーザが注視するシーンに関する情報との相関関係を学習したニューラルネットワークを利用して、前記ユーザが注視するシーンと相関関係のある情報を推定する、
請求項 8 に記載の人工知能情報処理装置。

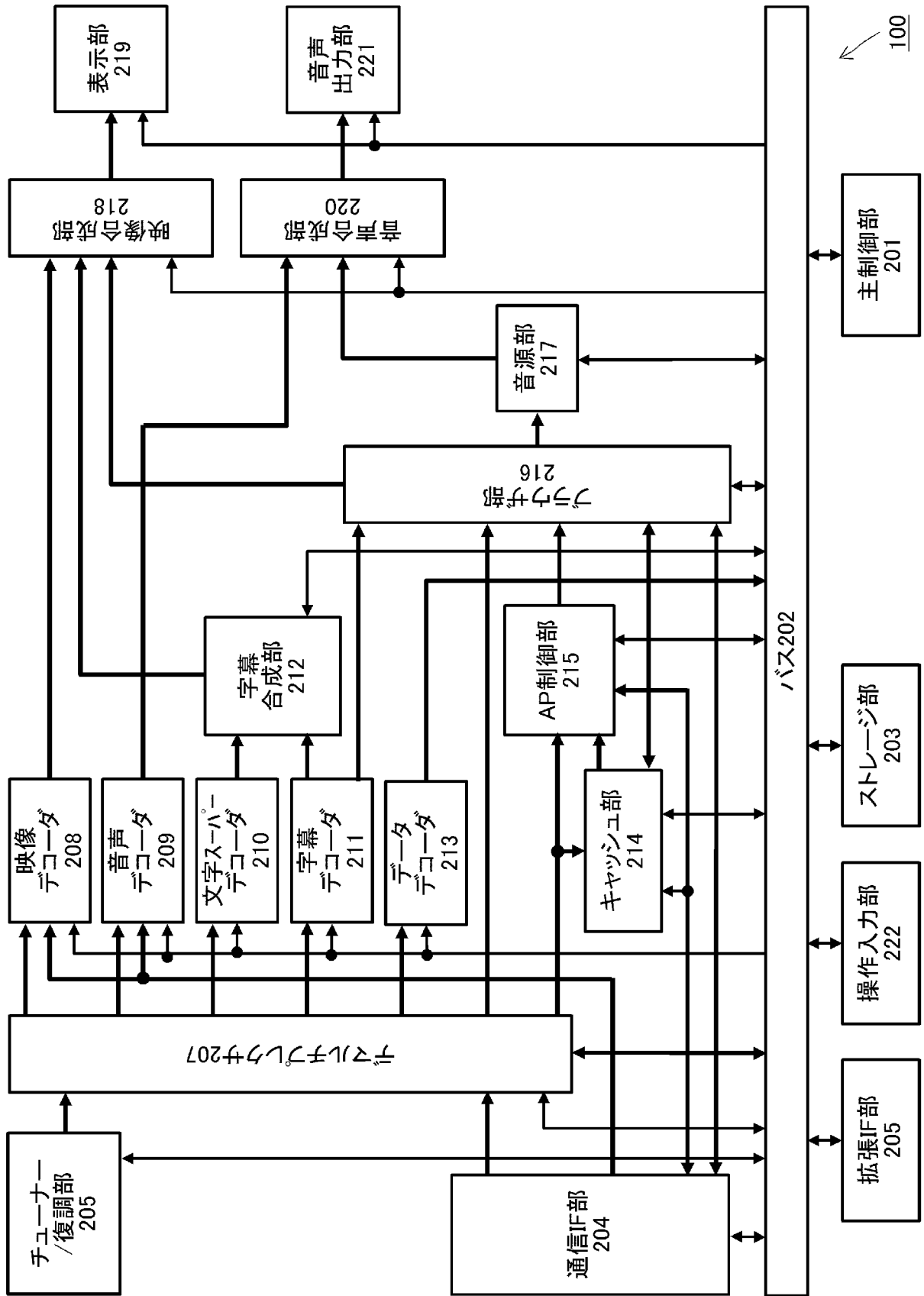
[請求項10] コンテンツを視聴しているユーザの注視度をセンサー情報に基づいて人工知能により推定する注視度推定ステップと、
前記注視度推定ステップにおける推定結果に基づいて、前記コンテンツ中で前記ユーザが注視するシーンの映像と前記コンテンツに関する情報を取得する取得ステップと、
前記ユーザが注視するシーンの映像と前記コンテンツに関する情報に基づいて、前記ユーザが注視しているシーンに関する情報を人工知

能により推定するシーン情報推定ステップと、
を有する人工知能情報処理方法。

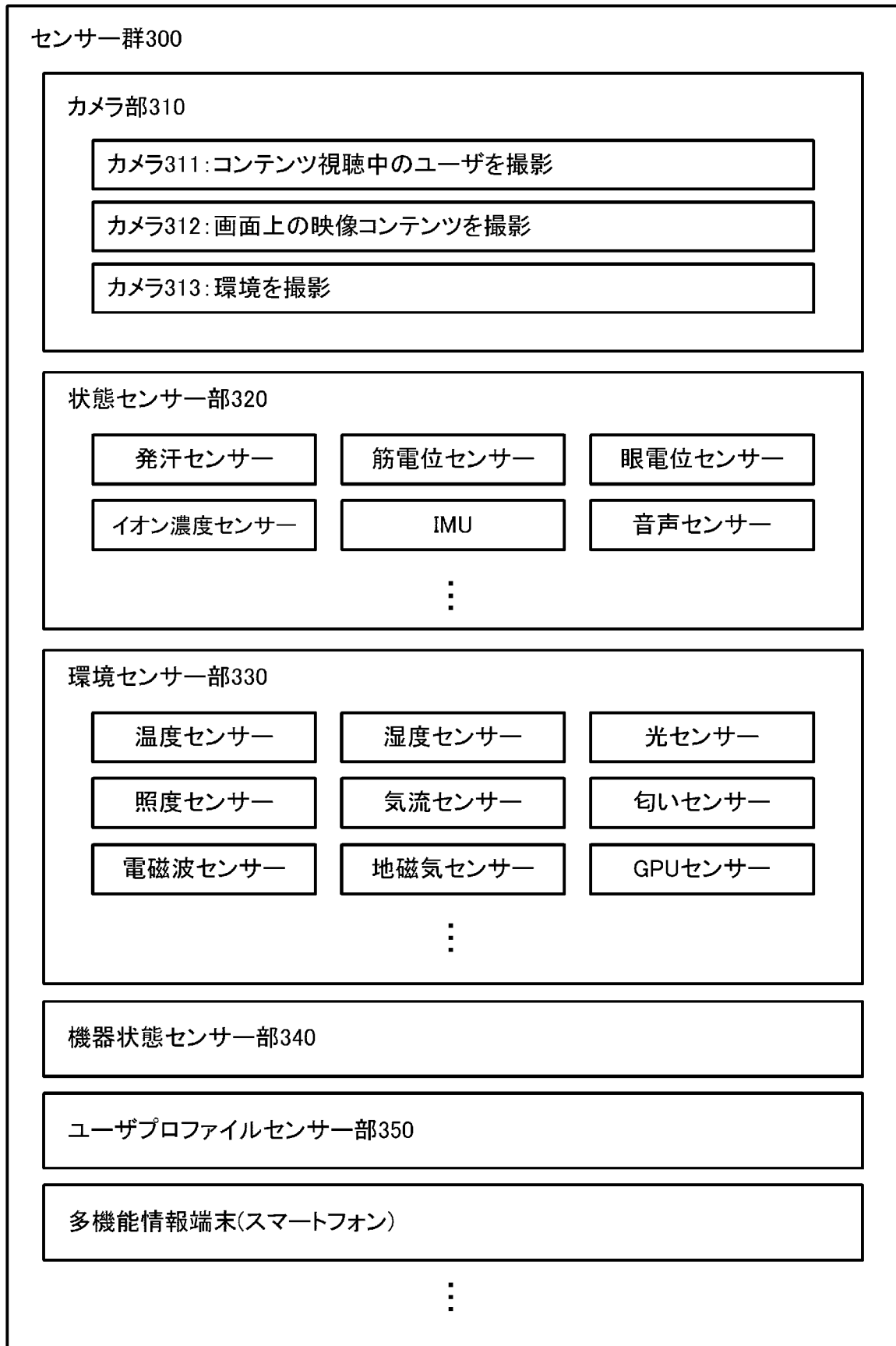
[図1]



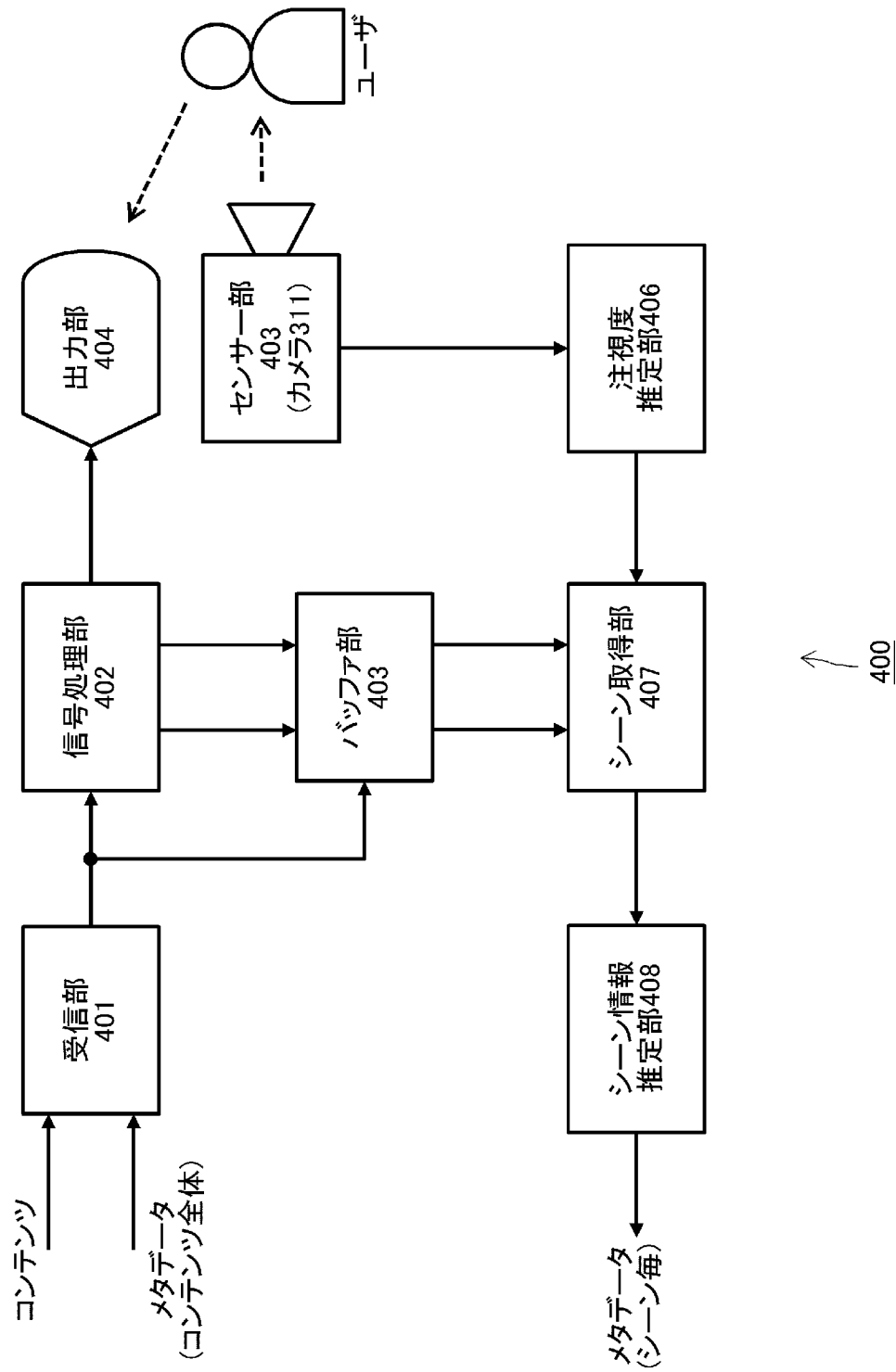
[図2]



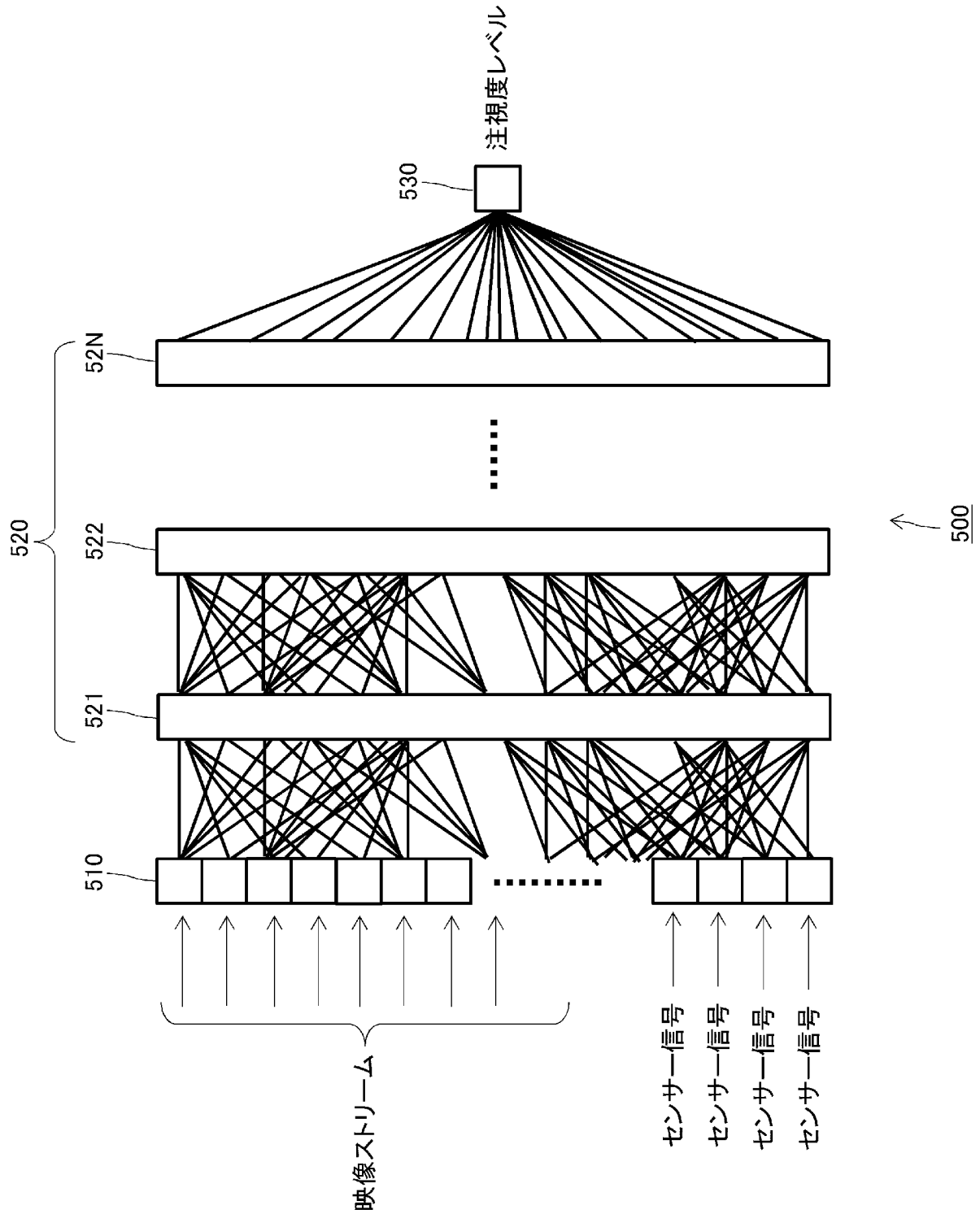
[図3]



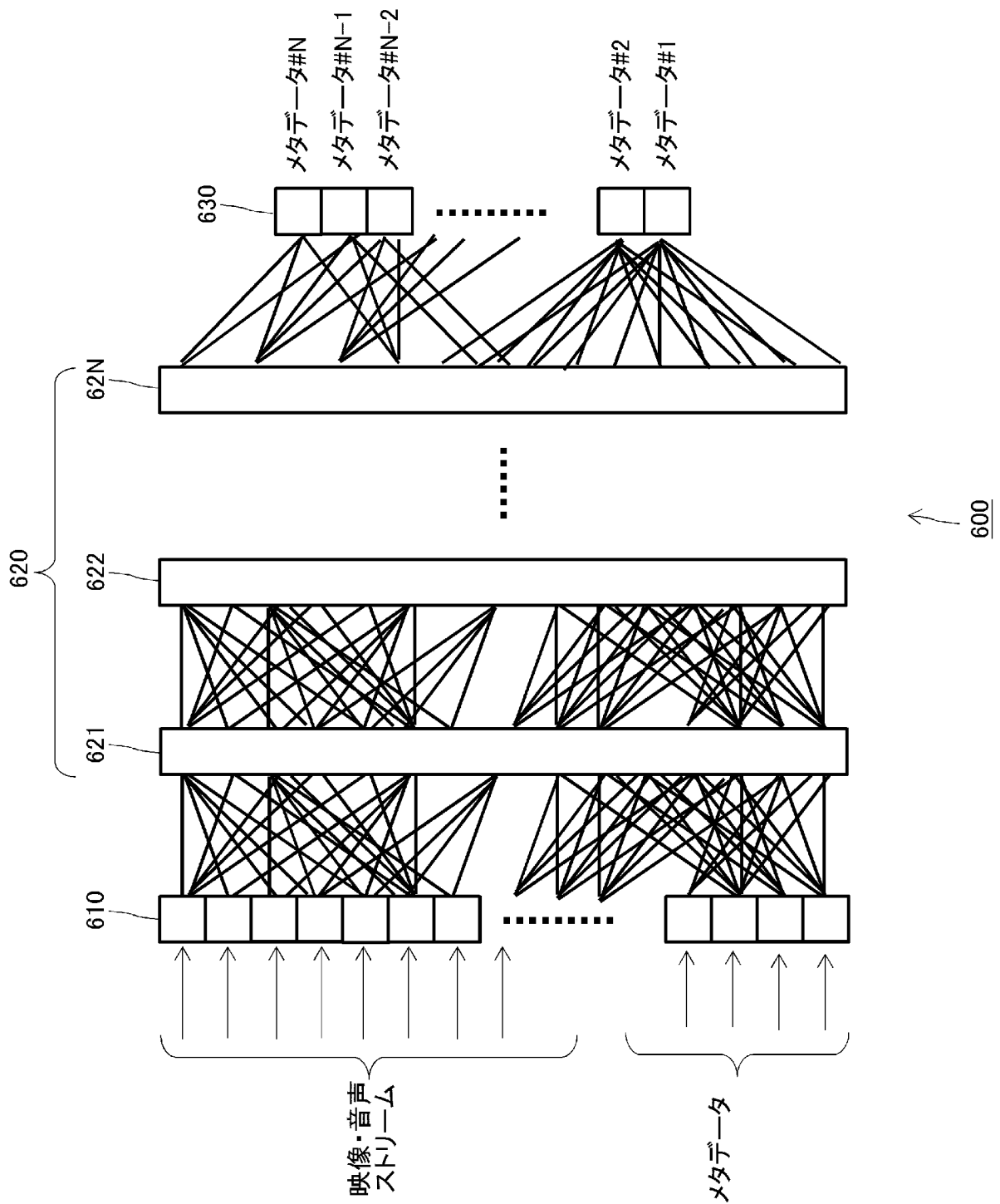
[図4]



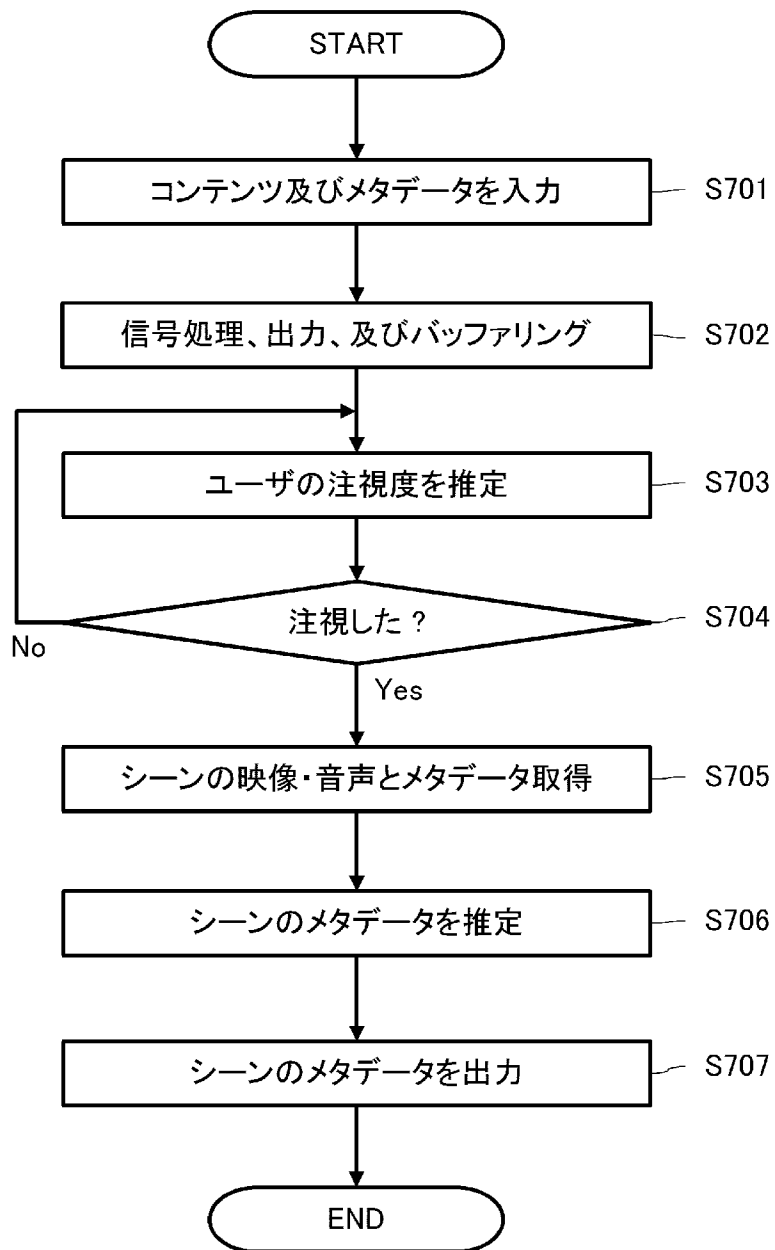
[図5]



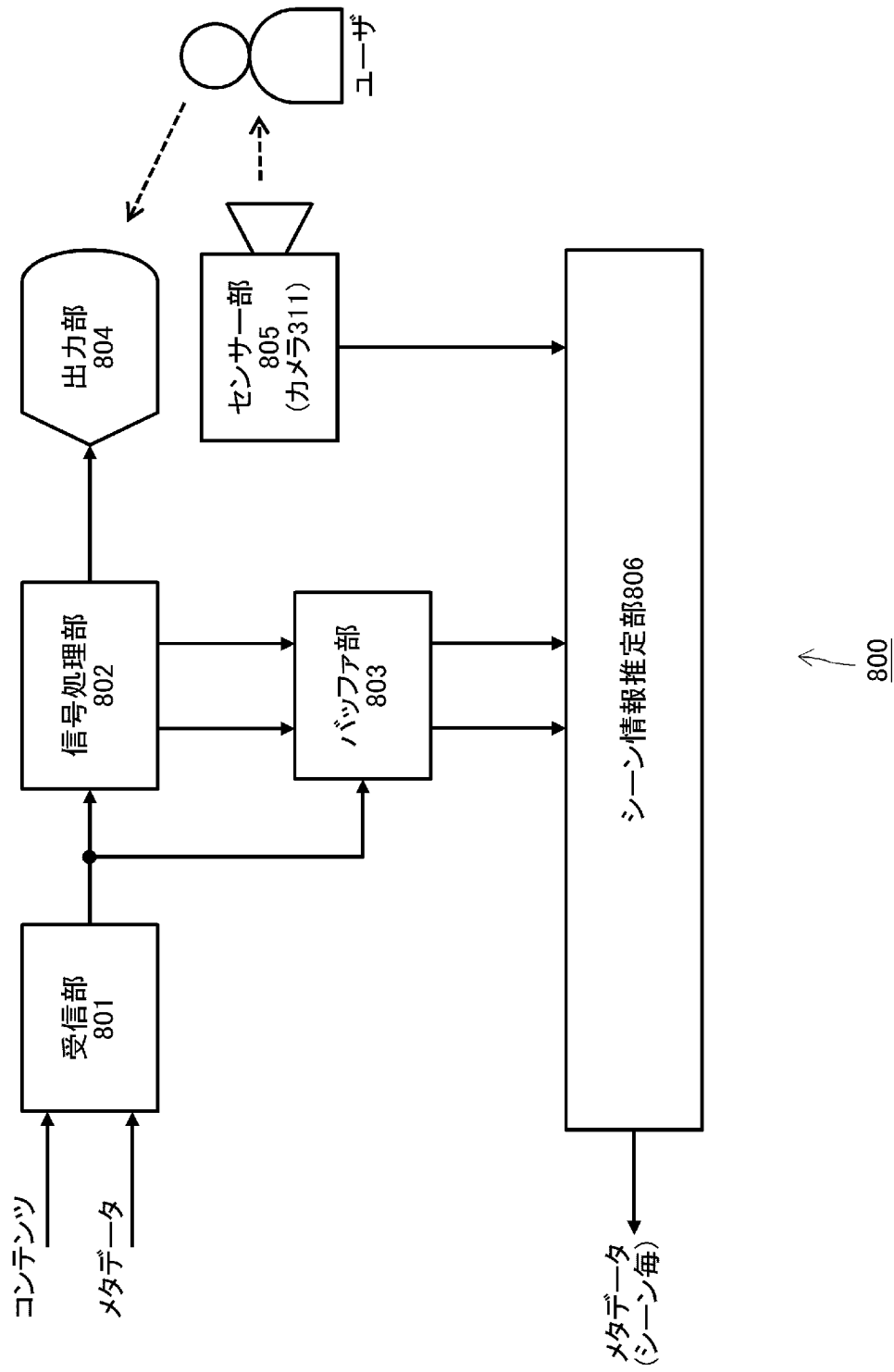
[図6]



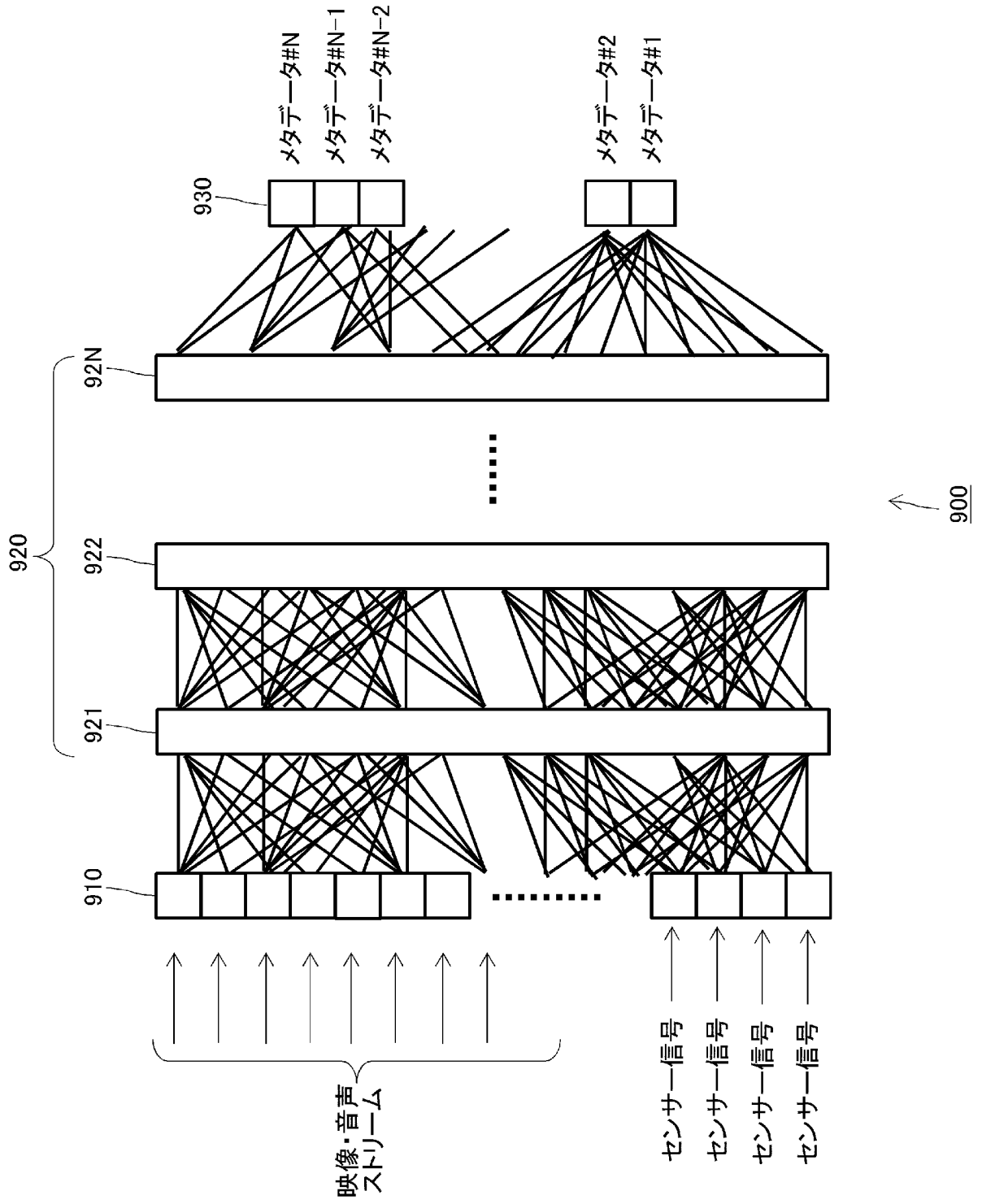
[図7]



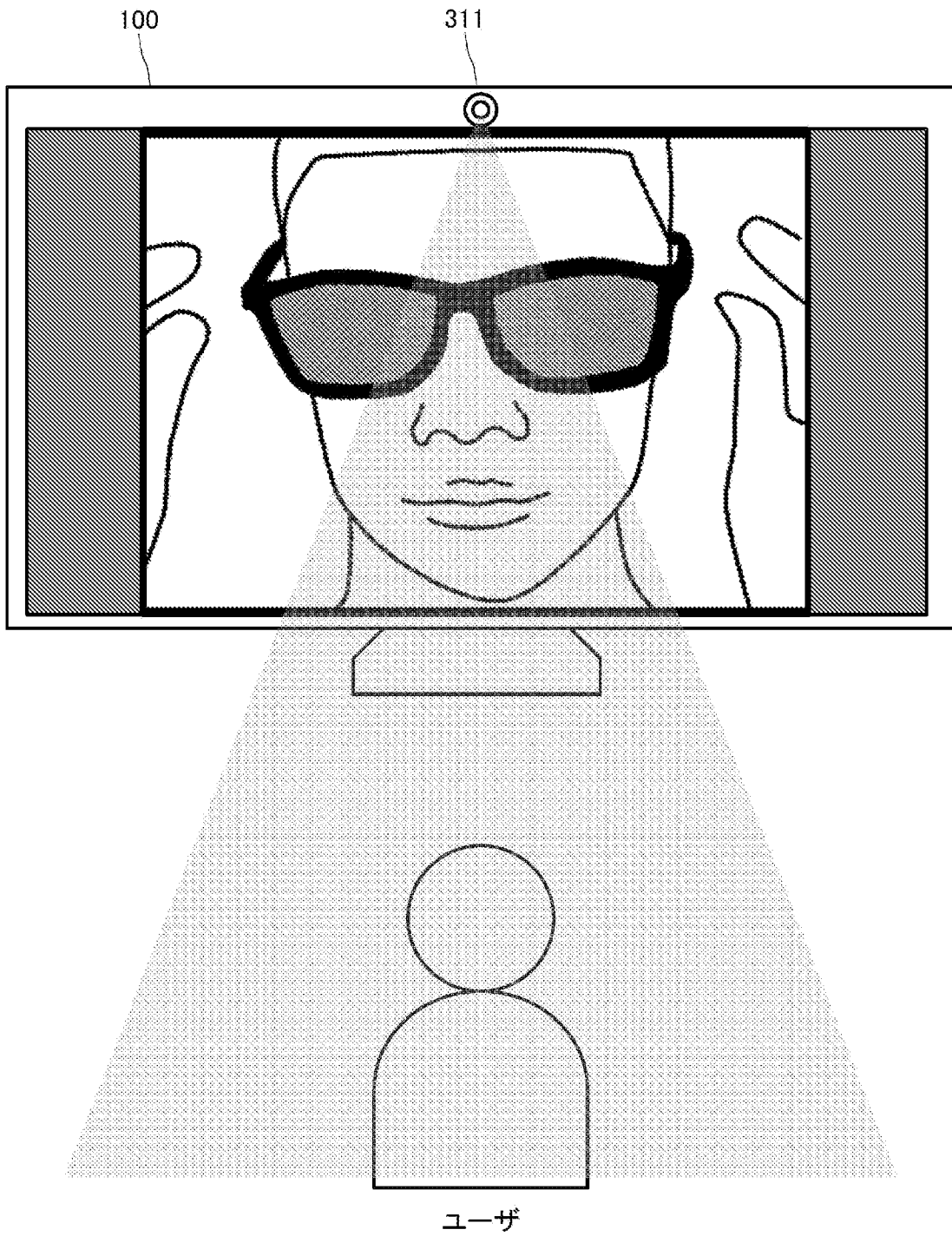
[図8]



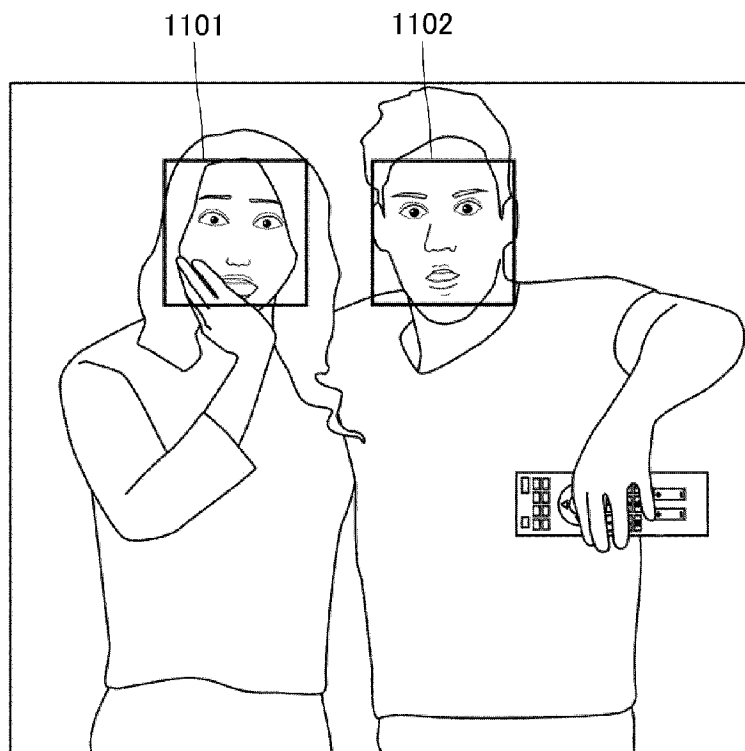
[図9]



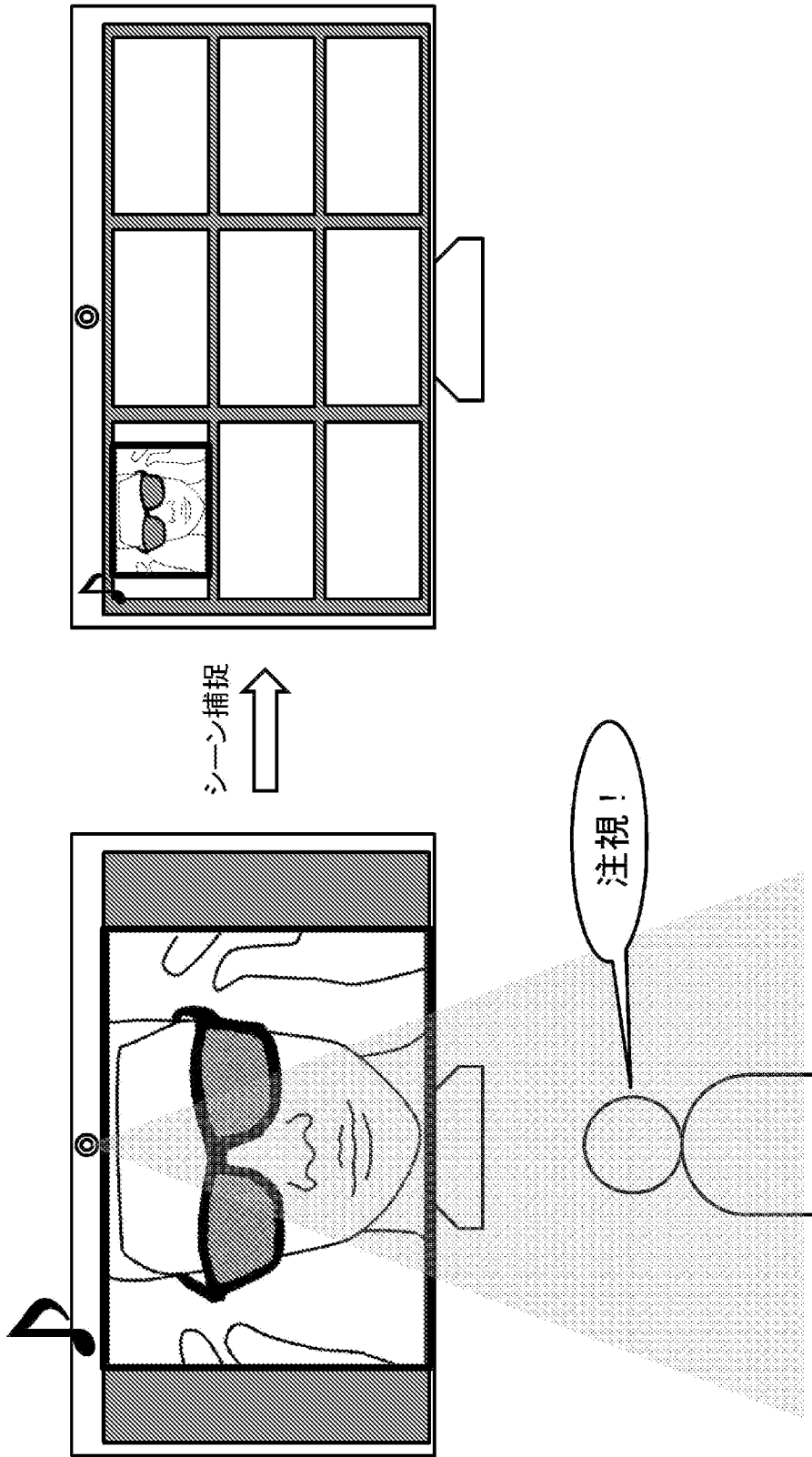
[図10]



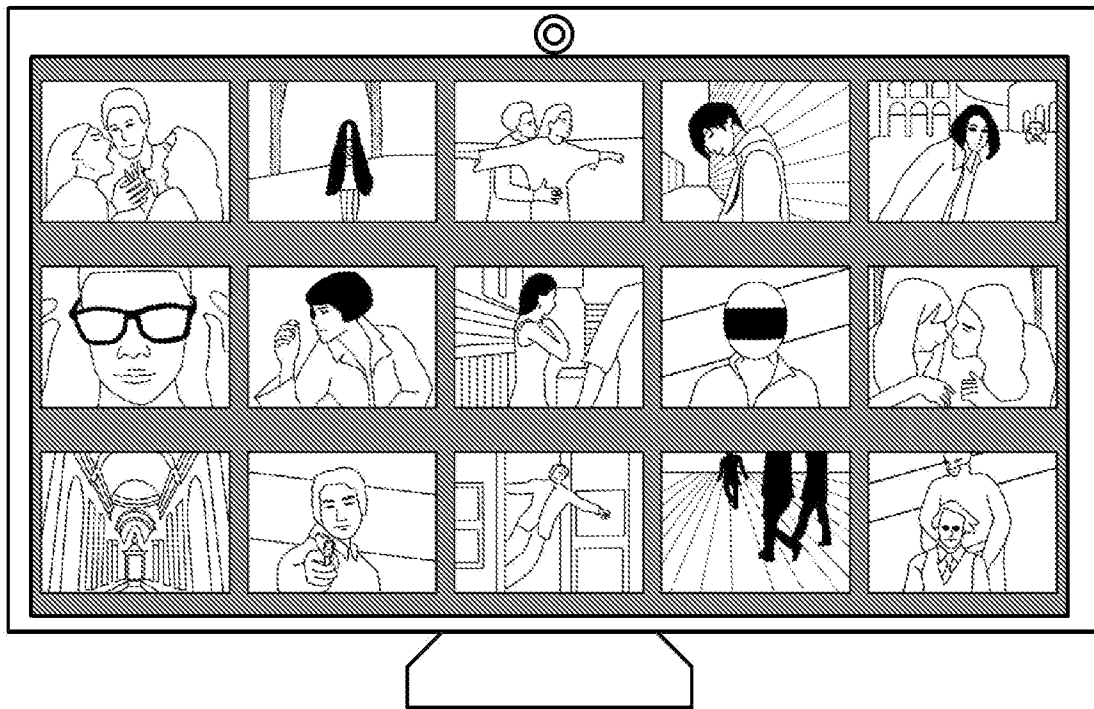
[図11]



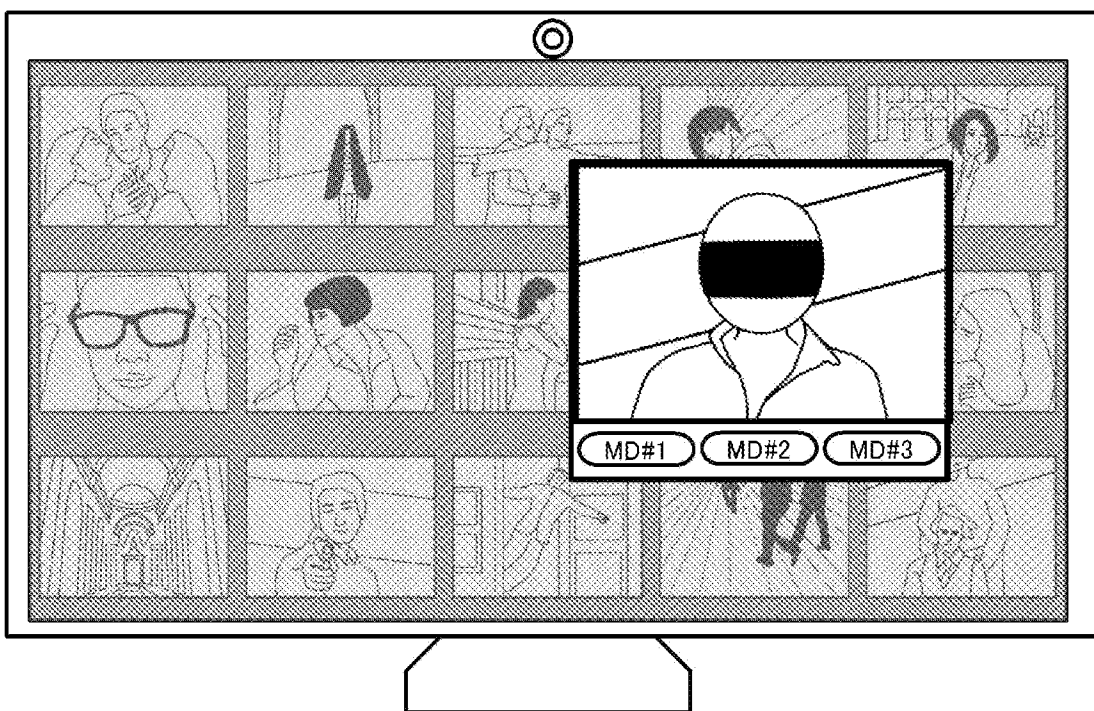
[図12]



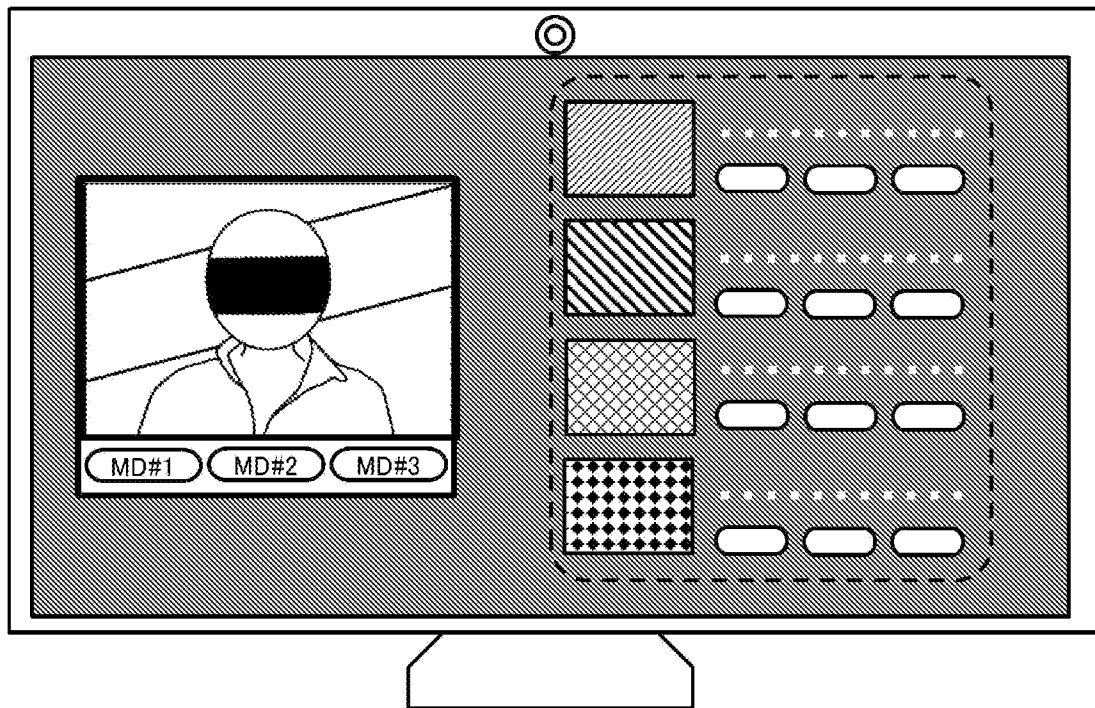
[図13]



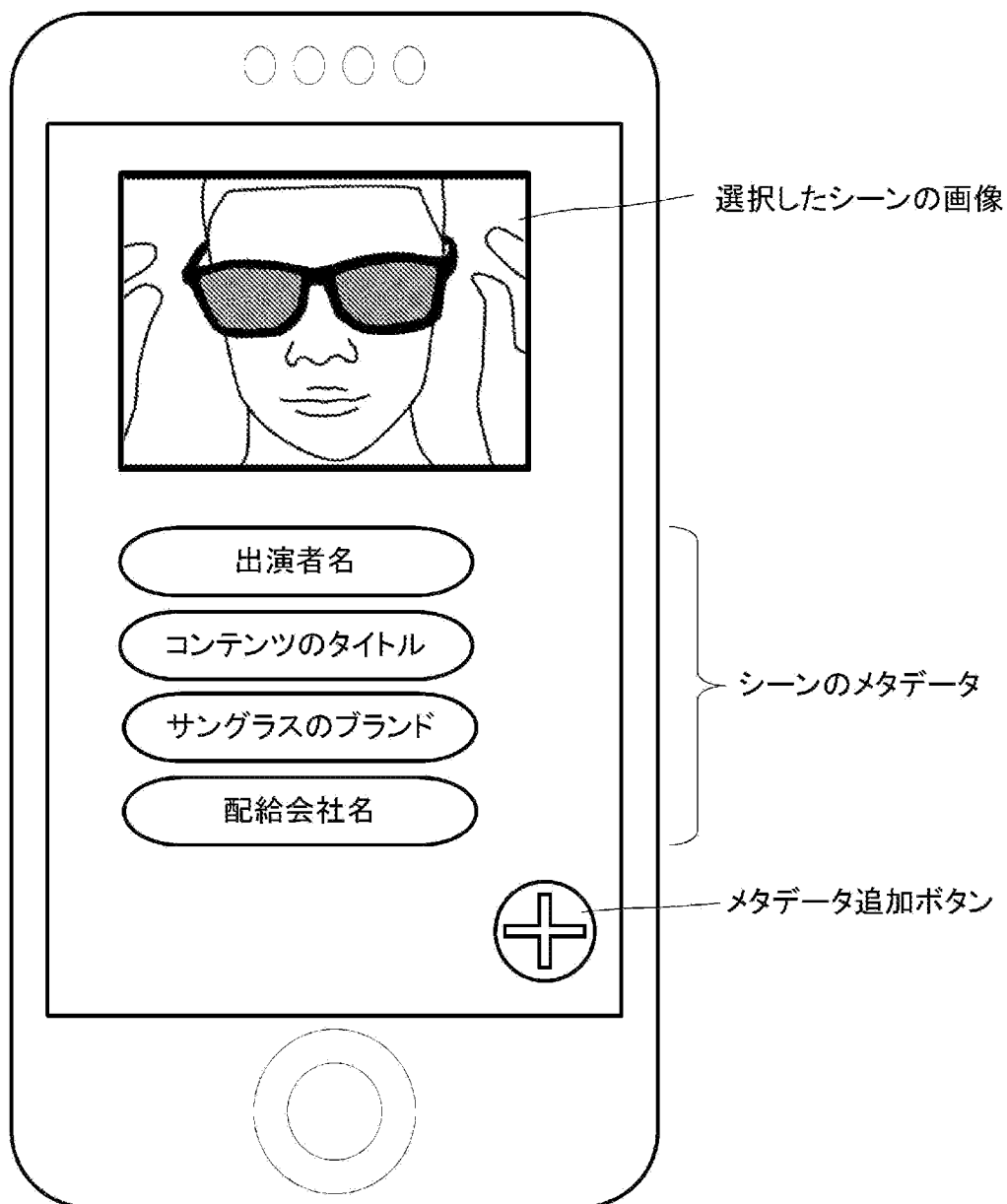
[図14]



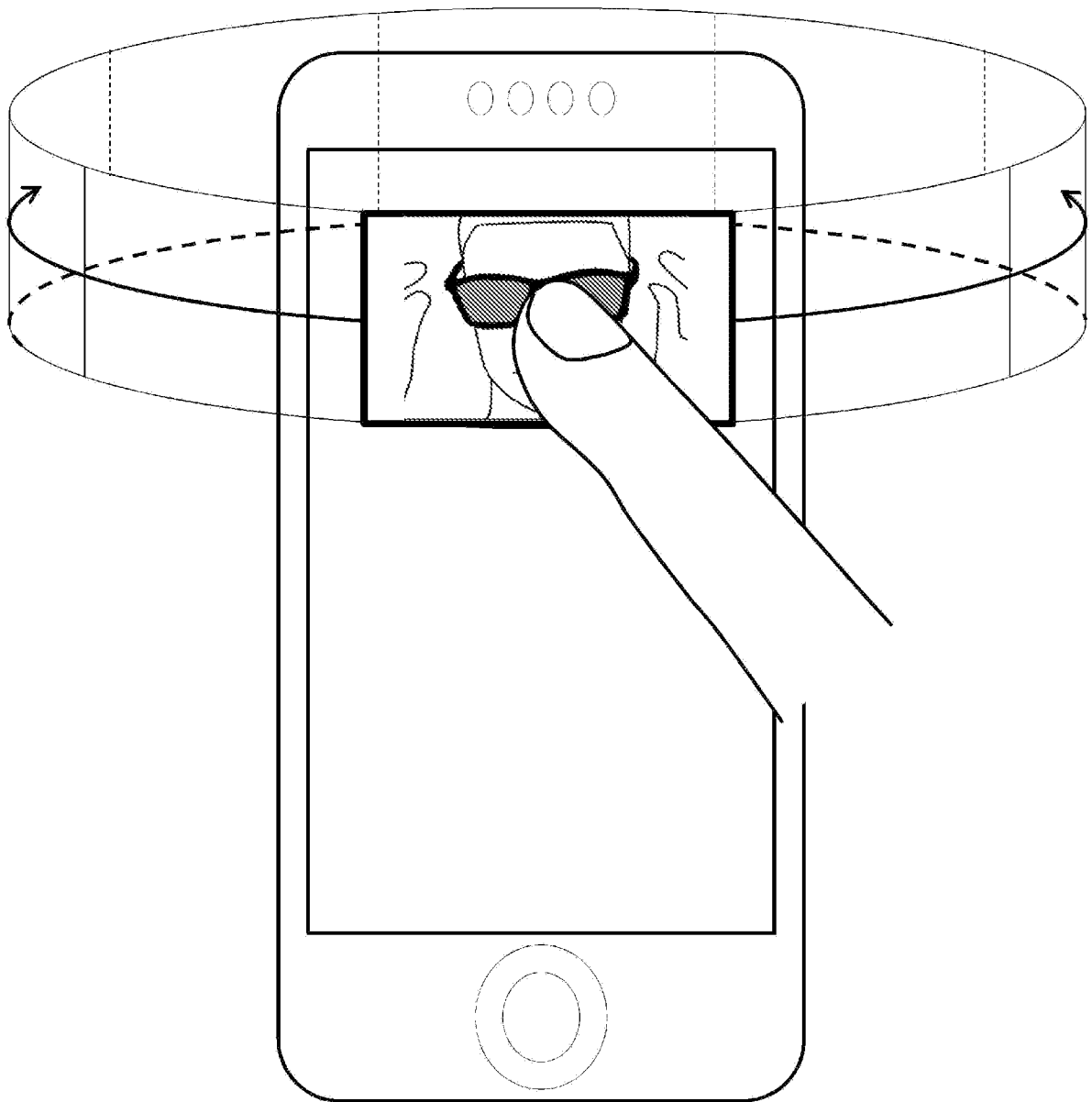
[図15]



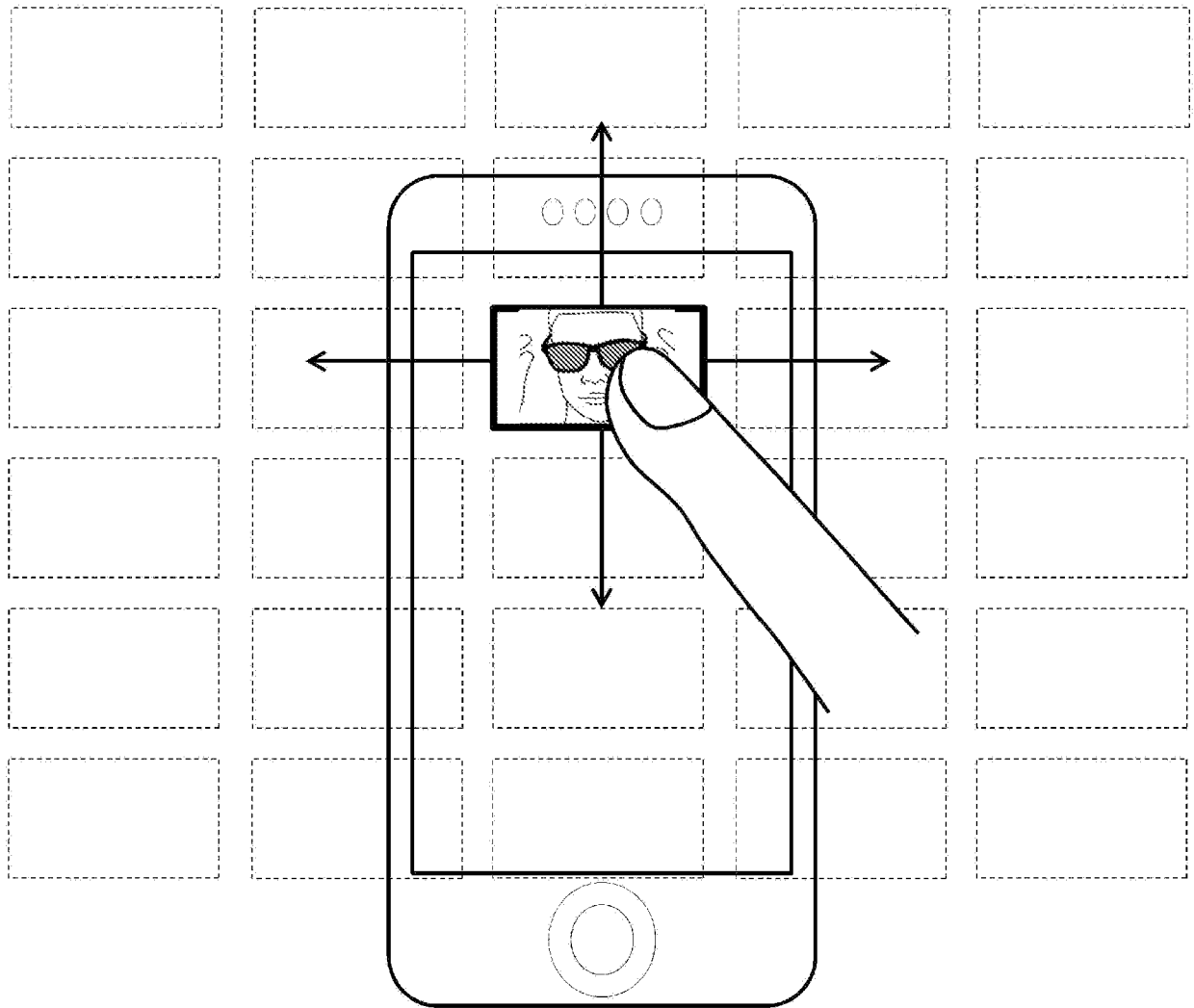
[図16]



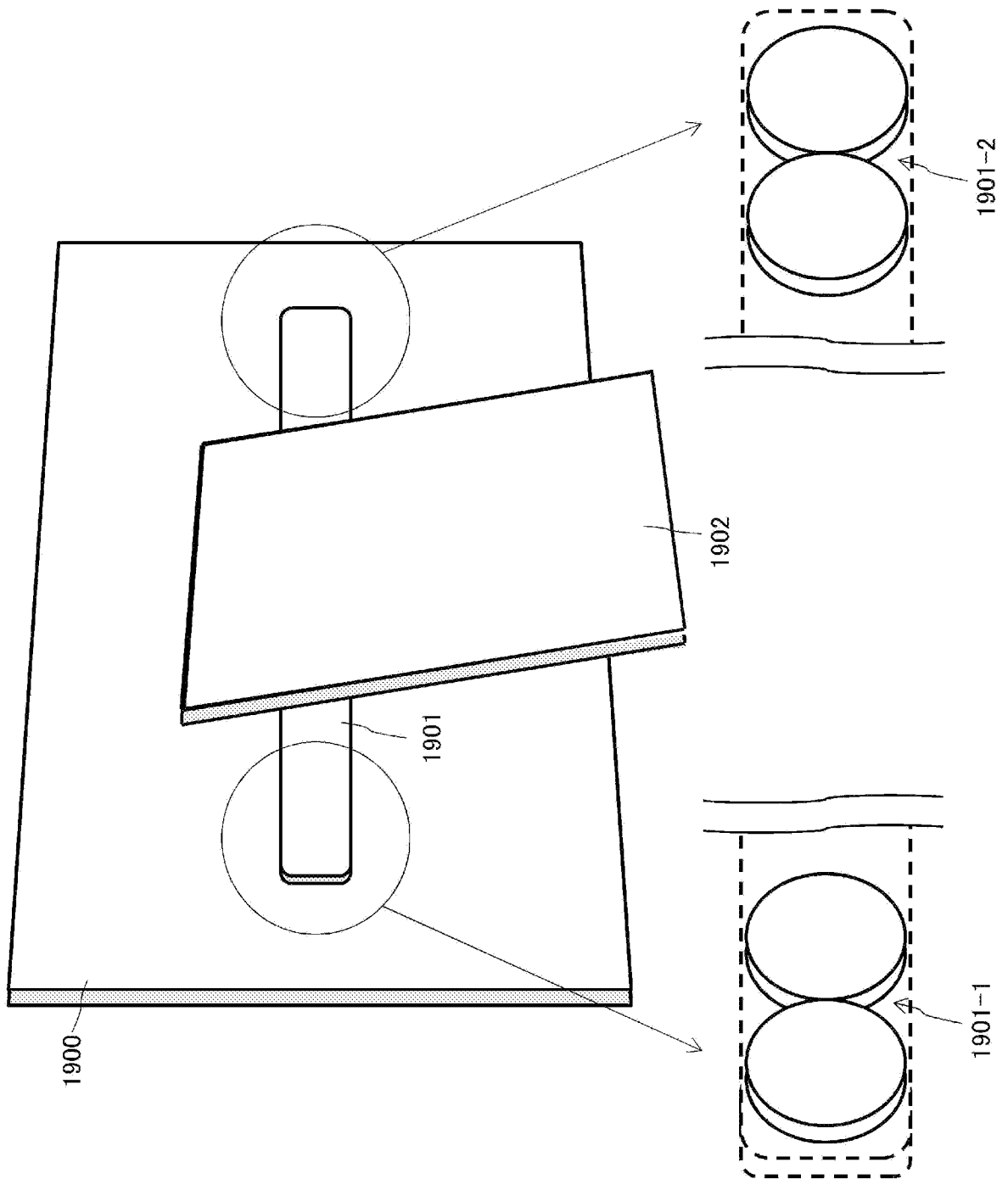
[図17]



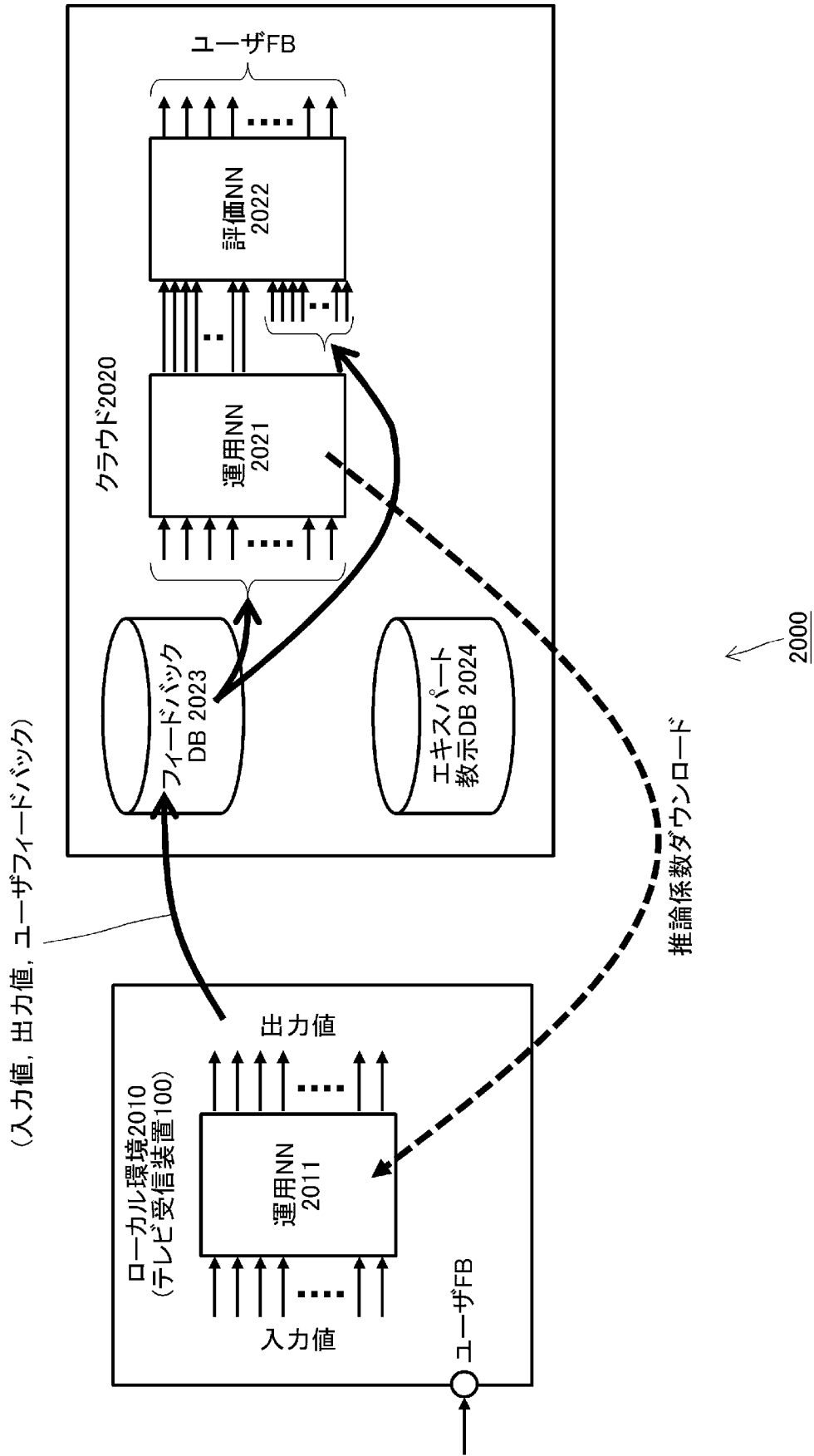
[図18]



[図19]



[図20]



INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2020/009957

A. CLASSIFICATION OF SUBJECT MATTER
 Int.Cl. G06F16/9035(2019.01) i, G06F16/9038(2019.01) i, G06F16/907(2019.01) i, H04N21/442(2011.01) i, H04N21/466(2011.01) i
 FI: H04N21/466, G06F16/9035, G06F16/9038, G06F16/907, H04N21/442
 According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)
 Int.Cl. G06F16/9035, G06F16/9038, G06F16/907, H04N21/442, H04N21/466

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Published examined utility model applications of Japan	1922-1996
Published unexamined utility model applications of Japan	1971-2020
Registered utility model specifications of Japan	1996-2020
Published registered utility model applications of Japan	1994-2020

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y A	JP 2005-142975 A (NIPPON TELEGRAPH AND TELEPHONE CORPORATION) 02.06.2005 (2005-06-02), paragraphs [0046]-[0096], fig. 1-9	8-10 1-7
Y A	JP 2018-205819 A (FUJITSU LIMITED) 27.12.2018 (2018-12-27), paragraphs [0011]-[0084], fig. 1-9	8-10 1-7
A	JP 2011-239158 A (JAPAN BROADCASTING CORPORATION) 24.11.2011 (2011-11-24), paragraphs [0036]-[0078], fig. 1-5	1-10

Further documents are listed in the continuation of Box C. See patent family annex.

* Special categories of cited documents:	“T” later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
“A” document defining the general state of the art which is not considered to be of particular relevance	“X” document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
“E” earlier application or patent but published on or after the international filing date	“Y” document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
“L” document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	“&” document member of the same patent family
“O” document referring to an oral disclosure, use, exhibition or other means	
“P” document published prior to the international filing date but later than the priority date claimed	

Date of the actual completion of the international search 13.05.2020	Date of mailing of the international search report 26.05.2020
---	--

Name and mailing address of the ISA/ Japan Patent Office 3-4-3, Kasumigaseki, Chiyoda-ku, Tokyo 100-8915, Japan	Authorized officer Telephone No.
--	---

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.
PCT/JP2020/009957

JP 2005-142975 A	02.06.2005	(Family: none)
JP 2018-205819 A	27.12.2018	(Family: none)
JP 2011-239158 A	24.11.2011	(Family: none)

A. 発明の属する分野の分類（国際特許分類（IPC）） G06F 16/9035(2019.01)i; G06F 16/9038(2019.01)i; G06F 16/907(2019.01)i; H04N 21/442(2011.01)i; H04N 21/466(2011.01)i FI: H04N21/466; G06F16/9035; G06F16/9038; G06F16/907; H04N21/442		
B. 調査を行った分野 調査を行った最小限資料（国際特許分類（IPC）） G06F16/9035; G06F16/9038; G06F16/907; H04N21/442; H04N21/466 最小限資料以外の資料で調査を行った分野に含まれるもの 日本国実用新案公報 1922 - 1996年 日本国公開実用新案公報 1971 - 2020年 日本国実用新案登録公報 1996 - 2020年 日本国登録実用新案公報 1994 - 2020年		
国際調査で使用した電子データベース（データベースの名称、調査に使用した用語）		
C. 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
Y A	JP 2005-142975 A（日本電信電話株式会社）02.06.2005（2005 - 06 - 02） [0046]-[0096], 図1-9	8-10 1-7
Y A	JP 2018-205819 A（富士通株式会社）27.12.2018（2018 - 12 - 27） [0011]-[0084], 図1-9	8-10 1-7
A	JP 2011-239158 A（日本放送協会）24.11.2011（2011 - 11 - 24） [0036]-[0078], 図1-5	1-10
<input type="checkbox"/> C欄の続きにも文献が列挙されている。 <input checked="" type="checkbox"/> パテントファミリーに関する別紙を参照。		
* 引用文献のカテゴリー “A” 特に関連のある文献ではなく、一般的な技術水準を示すもの “E” 国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの “L” 優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献（理由を付す） “O” 口頭による開示、使用、展示等に言及する文献 “P” 国際出願日前で、かつ優先権の主張の基礎となる出願の日の後に公表された文献	“T” 国際出願日又は優先日後に公表された文献であって出願と抵触するものではなく、発明の原理又は理論の理解のために引用するもの “X” 特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの “Y” 特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの “&” 同一パテントファミリー文献	
国際調査を完了した日 13.05.2020	国際調査報告の発送日 26.05.2020	
名称及びあて先 日本国特許庁(ISA/JP) 〒100-8915 日本国 東京都千代田区霞が関三丁目4番3号	権限のある職員（特許庁審査官） 松元 伸次 5C 9563 電話番号 03-3581-1101 内線 3541	

国際調査報告
特許ファミリーに関する情報

国際出願番号

PCT/JP2020/009957

引用文献	公表日	特許ファミリー文献	公表日
JP 2005-142975 A	02.06.2005	(ファミリーなし)	
JP 2018-205819 A	27.12.2018	(ファミリーなし)	
JP 2011-239158 A	24.11.2011	(ファミリーなし)	