

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
16 February 2012 (16.02.2012)

(10) International Publication Number
WO 2012/021541 A1

(51) International Patent Classification:
H04N 5/232 (2006.01) *H04N 13/00* (2006.01)
H04N 13/02 (2006.01)

(74) Agent: FULLER, Michael, L.; KNOBBE MARTENS
OLSON & BEAR, LLP, 2040 Main Street, 14th Floor,
Irvine, CA 92614 (US).

(21) International Application Number:
PCT/US2011/047126

(81) Designated States (unless otherwise indicated, for every
kind of national protection available): AE, AG, AL, AM,
AO, AT, AU, AZ, BA, BB, BG, BH, BR, BW, BY, BZ,
CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO,
DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT,
HN, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KN, KP,
KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD,
ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI,
NO, NZ, OM, PE, PG, PH, PL, PT, QA, RO, RS, RU,
SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM,
TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM,
ZW.

(22) International Filing Date:
9 August 2011 (09.08.2011)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:
61/371,979 9 August 2010 (09.08.2010) US
61/489,231 23 May 2011 (23.05.2011) US
13/205,481 8 August 2011 (08.08.2011) US

(71) Applicant (for all designated States except US): QUAL-
COMM INCORPORATED [US/US]; 5775 Morehouse
Drive, San Diego, CA 92121 (US).

(84) Designated States (unless otherwise indicated, for every
kind of regional protection available): ARIPO (BW, GH,
GM, KE, LR, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG,
ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, MD, RU, TJ,
TM), European (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, ML, MR, NE, SN, TD, TG).

(72) Inventors; and

(75) Inventors/Applicants (for US only): ATANASSOV,
Kalin, M. [US/US]; 5775 Morehouse Drive, San Diego,
CA 92121 (US). GOMA, Sergiu, R. [US/US]; 5775
Morehouse Drive, San Diego, CA 92121 (US). RA-
MACHANDRA, Vikas [IN/US]; 5775 Morehouse Drive,
San Diego, CA 92121 (US).

Published:

— with international search report (Art. 21(3))

(54) Title: AUTOFOCUS FOR STEREOSCOPIC CAMERA

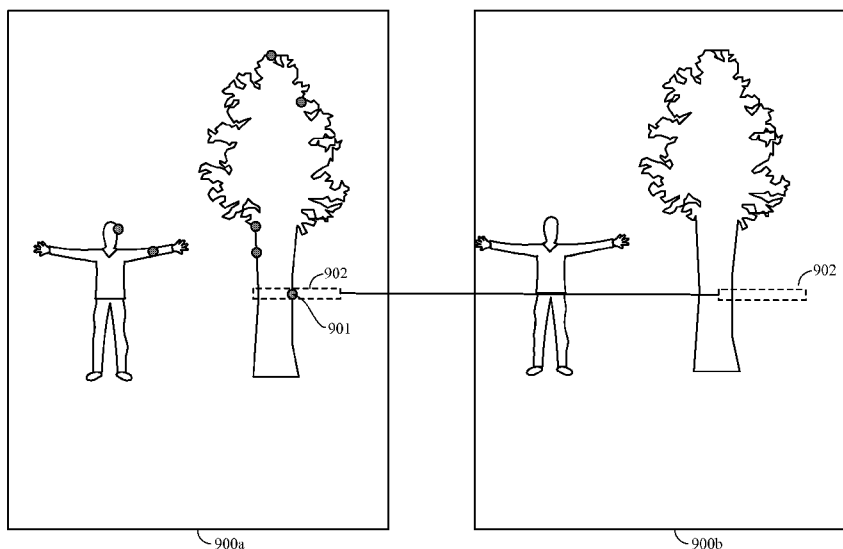


FIG. 9

(57) Abstract: Present embodiments contemplate systems, apparatus, and methods to determine an appropriate focal depth for a sensor based upon a pair of stereoscopic images. Particularly, certain of the embodiments contemplate determining keypoints for each image, identifying correlations between the keypoints, and deriving object distances from the correlations. These distances may then be used to select a proper focal depth for one or more sensors.

WO 2012/021541 A1

AUTOFOCUS FOR STEREOSCOPIC CAMERA

CROSS REFERENCE TO RELATED APPLICATIONS

[0001] The present application claims priority to U.S. Patent Application No. 13/205,481, filed on August 8, 2011, which claims the benefit of U.S. Provisional Patent Application Serial No. 61/371,979, entitled "INSTANTANEOUS AUTOFOCUS FROM STEREO IMAGES," filed August 9, 2010 and U.S. Provisional Patent Application Serial No. 61/489,231, entitled "AUTOFOCUS FOR STEREO IMAGES," filed on May 23, 2011, which applications are incorporated by reference herein.

TECHNICAL FIELD

[0002] The present embodiments relate to stereoscopic image capture, and in particular, to methods, apparatus and systems for determining an appropriate focal depth for a stereoscopic image capture device.

BACKGROUND

[0003] Stereopsis is the process by which the human brain interprets an object's depth based on the relative displacement of the object as seen from the left and right eyes. The stereoscopic effect may be artificially induced by taking first and second images of a scene from first and second laterally offset viewing positions and presenting the images separately to each of the left and right eyes. By capturing a succession of stereoscopic image pairs in time, the image pairs may be successively presented to the eyes to form a stereoscopic movie that appears to the user as having three-dimensions.

[0004] Two traditional cameras may be used to acquire each of the images of a stereoscopic image pair. A traditional camera may be properly focused using an autofocus procedure which captures a plurality of images at different focal depths. The focal depth corresponding to the highest frequency content is then used for subsequent image captures. Traditional movie cameras may use this method to autofocus during

video capture. However, the capture of frames will need to be periodically delayed while the autofocus functionality is performed.

[0005] While suitable for capturing 2D images with a single traditional camera, this autofocus technique may be unsuitable for stereoscopic image capture. In particular, the technique may disrupt the video stream and may be affected by camera movement, such as by the user's hand motions.

SUMMARY

[0006] Certain embodiments contemplate a method in an electronic device for determining a focal depth for an image sensor. The method may comprise: receiving a first image associated with a first viewpoint; receiving a second image associated with a second viewpoint; determining a first plurality of keypoints based on the first image; correlating keypoints from the first plurality of keypoints with positions in the second image; determining a plurality of disparities associated with each of the first plurality of keypoints; and determining a focal depth based upon the plurality of disparities, the position of the first viewpoint and the position of the second viewpoint.

[0007] In certain embodiments, the method may further comprise determining a second plurality of keypoints based on the second image. In some embodiments, correlating keypoints from the first plurality of keypoints with positions in the second image comprises correlating keypoints from the first plurality of keypoints with keypoints from the second plurality of keypoints. In some embodiments, correlating keypoints from the first plurality of keypoints with positions in the second image comprises iterating over pixels within a search range in the second image. In some embodiments, correlating keypoints from the first plurality of keypoints with positions in the second image comprises determining the mean square error between pixels in the first image and the second image. In some embodiments, determining a first plurality of keypoints based on the first image comprises determining Scale Invariant Feature Transform (SIFT) keypoints based on the first image. In some embodiments, determining a first plurality of keypoints based on the first image comprises sub-sampling the first image, applying a high-pass filter to the first image, calculating the power of the first image, and thresholding the first image. In some embodiments, correlating keypoints from the first plurality of keypoints with positions

in the second image occurs in realtime. In some embodiments, the electronic device comprises a mobile phone.

[0008] Certain embodiments contemplate a computer readable medium comprising instructions configured to cause a computer to perform the steps of: receiving a first image associated with a first viewpoint; receiving a second image associated with a second viewpoint; determining a first plurality of keypoints based on the first image; correlating keypoints from the first plurality of keypoints with positions in the second image; determining a plurality of disparities associated with each of the first plurality of keypoints; and determining a focal depth based upon the plurality of disparities, the position of the first viewpoint and the position of the second viewpoint.

[0009] In some embodiments, the instructions are also configured to cause the processor to determine a second plurality of keypoints based on the second image. In some embodiments, correlating keypoints from the first plurality of keypoints with positions in the second image comprises correlating keypoints from the first plurality of keypoints with keypoints from the second plurality of keypoints. In some embodiments, correlating keypoints from the first plurality of keypoints with positions in the second image comprises iterating over pixels within a search range in the second image. In some embodiments, correlating keypoints from the first plurality of keypoints with positions in the second image comprises determining the mean square error between pixels in the first image and the second image. In some embodiments, determining a first plurality of keypoints based on the first image comprises determining Scale Invariant Feature Transform (SIFT) keypoints based on the first image. In some embodiments, determining a first plurality of keypoints based on the first image comprises sub-sampling the first image, applying a high-pass filter to the first image, calculating the power of the first image, and thresholding the first image. In some embodiments, correlating keypoints from the first plurality of keypoints with positions in the second image occurs in realtime. In some embodiments, the computer is located in a mobile phone.

[0010] Certain embodiments contemplate a system for focusing a stereoscopic capture device. The system may comprise a first image sensor configured to generate a first image associated with a first viewpoint; a second image sensor configured to generate a second image associated with a second viewpoint; a feature generation module configured to determine a first plurality of keypoints based on the

first image; a keypoint correlation module configured to correlate keypoints from the first plurality of keypoints with positions in the second image; a disparity determination module configured to determine a plurality of disparities associated with each of the first plurality of keypoints; and a depth determination module configured to determine a focal depth based upon the plurality of disparities, the position of the first viewpoint and the position of the second viewpoint.

[0011] In some embodiments, the feature generation module may be configured to determine a second plurality of keypoints based on the second image. In some embodiments, the software module configured to correlate keypoints is configured to correlate keypoints from the first plurality of keypoints with keypoints from the second plurality of keypoints. In some embodiments, the software module configured to correlate keypoints is configured to iterate over pixels within a search range in the second image. In some embodiments, correlating keypoints from the first plurality of keypoints with positions in the second image comprises determining the mean square error between pixels in the first image and the second image.

[0012] In some embodiments, the feature generation module is configured to determine Scale Invariant Feature Transform (SIFT) keypoints based on the first image. In some embodiments, the feature generation module is configured to sub-sample the first image, apply a high-pass filter to the first image, calculate the power of the first image, and threshold the first image.

[0013] In some embodiments, the software module configured to correlate keypoints correlates keypoints from the first plurality of keypoints with positions in the second image in realtime. In some embodiments, the stereoscopic capture device is located on a mobile phone. In some embodiments, the software module configured to determine a focal depth comprises a disparity histogram.

[0014] Certain embodiments contemplate a system for focusing a stereoscopic capture device, the system comprising: means for receiving a first image associated with a first viewpoint; means for receiving a second image associated with a second viewpoint; means for determining a first plurality of keypoints based on the first image; means for correlating keypoints from the first plurality of keypoints with positions in the second image; means for determining a plurality of disparities associated with each of the first plurality of keypoints; and means for determining a

focal depth based upon the plurality of disparities, the position of the first viewpoint and the position of the second viewpoint.

[0015] In some embodiments the means for receiving a first image comprises a first sensor, the means for receiving a second image comprises a second sensor, the means for determining a first plurality of keypoints comprises a feature generation module, the means for correlating comprises a keypoint correlation module, the means for determining a plurality of disparities comprises a disparity determination module, and the means for determining a focal depth comprises a depth determination module. In some embodiments, the means for determining a first plurality of keypoints is configured to determine a second plurality of keypoints based on the second image. In some embodiments, the means for correlating keypoints from the first plurality of keypoints with positions in the second image is configured to correlate keypoints from the first plurality of keypoints with keypoints from the second plurality of keypoints. In some embodiments, the means for correlating keypoints from the first plurality of keypoints with positions in the second image is configured to iterate over pixels within a search range in the second image. In some embodiments, the means for correlating keypoints from the first plurality of keypoints with positions in the second image is configured to determine the mean square error between pixels in the first image and the second image. In some embodiments, the means for determining a first plurality of keypoints is configured to determine Scale Invariant Feature Transform (SIFT) keypoints based on the first image. In some embodiments, the means for determining a first plurality of keypoints is configured to sub-sample the first image, apply a high-pass filter to the first image, calculate the power of the first image, and threshold the first image. In some embodiments, the means for correlating keypoints correlates the keypoints from the first plurality of keypoints with positions in the second image in realtime. In some embodiments, the stereoscopic capture device is located on a mobile phone.

BRIEF DESCRIPTION OF THE DRAWINGS

[0016] The disclosed aspects will hereinafter be described in conjunction with the appended drawings, provided to illustrate and not to limit the disclosed aspects, wherein like designations denote like elements.

[0017] FIG. 1 is a generalized diagram depicting one possible mobile device comprising a sensor arrangement facilitating the capture of stereoscopic images.

[0018] FIG. 2 is a block diagram of certain of the components in a mobile device, such as the mobile device of Figure 1.

[0019] FIG. 3 depicts the capturing of an object at a first and second position using a stereo pair of capture devices.

[0020] FIG. 4 is a graph depicting the relationship between object distance and pixel disparity for a particular camera arrangement.

[0021] FIG. 5A is a diagram depicting a top-down view of an arbitrary scene and two image capture sensors positioned so as to achieve a stereoscopic effect.

[0022] FIG. 5B depicts one of the pair of stereoscopic images taken of the scene in Figure 5A with the magnitude and direction of object disparities in the scene overlaid.

[0023] FIG. 6 depicts the graph of Figure 4, but with an object disparity histogram and a corresponding object depth histogram overlaid.

[0024] FIG. 7 depicts a flow diagram for the process by which certain embodiments determine a new focal depth.

[0025] FIG. 8 depicts a flow diagram for the process by which certain of the embodiments determine keypoints.

[0026] FIG. 9 depicts a stereoscopic image pair and a region in which keypoints are correlated between each of the images.

DETAILED DESCRIPTION

[0027] Embodiments relate to systems and methods of determining or setting configuration data in a stereoscopic camera. In one embodiment, the configuration data relates to the proper focal length of the two lenses of the stereoscopic camera. In one embodiment, a first camera receives a first image from a scene and a second camera receives a second image of the same scene. A set of keypoints are determined from analysis of the first image. The keypoints, can be, for example, Keypoints may comprise any data structure which can be consistently replicated from a portion of an image and thereby permit unique identification of the image portion. In some embodiments, a keypoint may comprise a plurality of pixels corresponding to a portion of an image. The keypoint may be associated with a position in the image. After determining a keypoint in the first image, the system looks for a similar position

in the second image. Once the similar position in the second image is identified, the system calculates the difference between the keypoints in the first image, and the corresponding position in the second image. This allows the system to determine the focal depth of the scene by knowing the disparity between the same keypoint positions in both frames, along with the positions of the stereoscopic lenses.

[0028] Present embodiments contemplate systems, apparatus, and methods to determine an appropriate focal depth for a sensor based upon at least a pair of stereoscopic images. Particularly, certain of the embodiments contemplate determining keypoints for each image, identifying correlations between the keypoints, and deriving object distances from the correlations. One skilled in the art will recognize that these embodiments may be implemented in hardware, software, firmware, or any combination thereof. The stereoscopic system may be implemented on a wide range of electronic devices, including mobile wireless communication devices, personal digital assistants (PDAs), laptop computers, desktop computers, digital cameras, digital recording devices, and the like.

[0029] Fig. 1 depicts a mobile device 100 comprising a sensor arrangement facilitating the capture of stereoscopic images, or other means for receiving an image. Such a device may be a mobile phone, personal digital assistant, gaming device, or the like. The device 100 may comprise a first sensor 101a and a second sensor 101b separated by a distance d . The device may also comprise user input controls 102 and a display 103. In some embodiments, the sensors 101a and 101b may be situated such that they are horizontally, but not vertically, offset when the users holds device 100 so as to capture a stereoscopic picture or movie.

[0030] Although this particular device depicts two sensors 101a and 101b one skilled in the art may readily conceive of a stereoscopic image capture device which comprises more or less than two image sensors. For example, a device with only a single sensor may operate in combination with a series of lenses or reflecting surfaces to acquire two images at the positions of sensors 101a and 101b in rapid succession. This arrangement would likewise be able to acquire a stereoscopic image pair for use with the methods described below and the single sensor could be focused accordingly. Thus, the methods and systems discussed in this application will be applicable to any system which acquires two images from a first and second viewpoint, so long as those viewpoints facilitate a stereoscopic depiction of the image scene. Thus, reference to a

pair of image sensors should not be considered to exclude the possibility of a single image sensor receiving images from two viewpoints.

[0031] Fig. 2 is a block diagram of certain of the components in a mobile device, such as the mobile device 100 depicted in Figure 1. Sensor 101a receives a first image of the stereoscopic image pair and sensor 101b receives a second image of the stereoscopic image pair. In some embodiments, the sensors may receive the images simultaneously. The device may comprise a video front end 102 and memory 103. Video front end 102 may process incoming raw image data from sensors 101a and 101b and store the data in memory 103. Memory 103 may also comprise various applications and software drivers for the mobile device 100. For example, a display driver module 104 may be in communication with the display 103. A user input module 106 may similarly be in communication with a user interface 102. A wireless communication driver module 107 may be in communication with wireless communications hardware 112.

[0032] The memory may also be in communication with a General Processor 113. The General Processor 113 may comprise sub-processing units, or subprocessors, such as an Advanced RISC Machine (ARM), digital signal processor (DSP), or graphical processing unit (GPU). These processors may communicate with local memory 114 when handling various operations.

[0033] Certain of the present embodiments contemplate the addition of a "Focal Depth Analysis Module" 115a, 115b to the system architecture. In some embodiments, the module may take the form of a dedicated processor 115a, or a portion of a processor located on the general processor. In some embodiments the module may comprise software code 115b stored in a computer readable medium such as memory 103. Some embodiments may place portions of the module at a dedicated processor 115a and memory 115b as firmware, or as a software-hardware combination. In some embodiments, the module may reside at any location in Figure 2 which permits access to a feature generation system, such as a SIFT feature generation system, and to sensors 101a and 101b. Thus, the module may take advantage of preexisting hardware or software configured for feature generation and/or detection. One skilled in the art will recognize that the embodiments described below could be implemented using a subprocessor on the General Processor 113, or could be stored as a separate application in memory 103. In some embodiments the SIFT feature generation system may be

found in software, whereas in other embodiments the SIFT feature generation system may be found in hardware.

[0034] Certain of the present embodiments provide auto-focus functionality which takes advantage of geometric properties of stereoscopic image capture. Fig. 3 depicts, via a top-down view, the stereoscopic image capture of an object 304 at a first position 300a and second position 300b using a stereoscopic camera arrangement. A first image capture device 301a may be located at a first position laterally separated from a second capture device 301b located at a second position. The first capture device 301a may capture a first image of the scene from the first position and the second capture device 301b may capture a second image of the scene from the second position. The first and second images will accordingly be associated with first and second viewpoints of the scene based on the positions and orientations of capture devices 301a and 301b. Object 304 may appear in both images. In some embodiments, capture device 301a and capture device 301b may be the same as sensors 101a and 101b of Figure 1 respectively. Capture devices 301a, 301b may be calibrated to have no vertical disparity and to possess fairly close focal distances.

[0035] The center of device 301a's viewpoint passes along line 302a. Similarly, the center of device 301b's viewpoint passes along line 302b. These two centerlines intersect at the position 303. As mentioned, the object 304 appears in each of the first and second images. With regard to position 300a, however, the object 304 appears to the right of centerline 302a, by an amount 305a and to the left of centerline 302b by an amount 305b. Conversely, in position 300b, the object 304 appears to the left of centerline 302a by an amount 306a, and to the right of centerline 302b, by an amount 306b. In this manner, the relative position of the object in the z-direction is reflected by the relative displacement in each of the left and right images.

[0036] Object disparity may be defined as the difference between an object's position in the first image as compared to the object's position in the second image. Where there is no vertical disparity between the capture devices, the disparity may comprise only the lateral offset from the position in one image to the position in another. One may arbitrarily take the disparity as the difference between the left and right, or right and left images. For the purposes of this description, the disparity is defined as the position of the object in the image from sensor 301b minus the position of the object in the image from sensor 301a (with the x-direction positive as indicated in

Figure 3). Thus, negative disparity results from the depiction of object 304 in the position 300a and positive disparity results from the depiction of object 304 in the position 300b.

[0037] With knowledge of the sensor positions and relative orientations one can construct a graph of the relationship between the observed disparity and an object's distance, or depth, from the camera arrangement. Figure 4, for example, is a graph of this relationship for one particular sensor arrangement. As the disparity 401 increases, the object distance 402 increases as well. An initial negative disparity may be present for objects very near the camera arrangement, i.e. those having little depth 402 in the z-direction of Figure 3. As the object is moved further from the camera arrangement (i.e. the depth increases) the disparity becomes increasingly positive, beginning to plateau for objects at a considerable distance. One may recognize that the chart of Figure 4 may depend upon the angles at which sensors 301a, 301b are oriented. Similarly, though the sensors may be parallel to one another as in Figures 1 and 3, displacement in the z and y directions between the sensors may also result in modifications to the graph. Such a graph may be stored in memory on the device, or in a similar storage structure for quick reference.

[0038] Figure 5a is a top-down view of a scene comprising several objects. Again, image capture devices 301a and 301b may be used to acquire images 501a and 501b respectively. Objects 502-504, are located at various depths within the scene. Consequently, disparities between the object positions in images 501a and 501b will be observed. Figure 5B depicts the image 501a with certain of the disparity magnitudes and directions indicated at the pixel positions for which they occur. For example, a plurality of positive disparities 510 emerge for the distant object 504, and a plurality of negative disparities 511 appear for the closer object 502. With reference to the graph of Figure 4, an automated system may determine the depth associated with each disparity. For example, as shown in Figure 6, disparities of the same magnitude have been accumulated and plotted to form the disparity histogram 601. The corresponding depths 602 may be derived from the relationship of the sensors to generate depth histogram 602. The depth histogram 602 would suggest the presence of one or more objects in the region of each maximum at the indicated depth.

[0039] An autofocus operation comprises the determination of the proper focal depth for one or both of the sensors 101a, 101b. In some embodiments, the proper

focal depth may be determined by taking the mean, median, or similar statistic of object depth histogram (or the object disparity histogram in conjunction with a graph such as Figure 4). The median statistic provides some robustness to outlying value while a special order statistic filter may be used to accommodate a particular application. The statistic selected may depend upon the relative weights to be given very distant and very close objects. For example, focal quality may be roughly the same through one range of depths, but vary dramatically in a second range. These variations are discussed in greater detail below.

[0040] While one could generate the disparity histogram 601 by determining the disparity of every pixel for every object found in each of the first and second images, this may be computationally expensive and impractical on a mobile device. Not only would the correlation of every pixel require iterating through a substantial number of pixels, but each image may comprise multiple pixels of the same value, making identification of an individual pixel and its correlation to an object in each image difficult.

[0041] In lieu of analyzing every pixel, certain of the present embodiments contemplate creating a “sparse” disparity map or “sparse” corresponding depth map of the image contents. In certain embodiments, keypoints may be determined in each of the images and the disparities between keypoints, or between keypoints and pixels, rather than between all or most of the pixels in the images, may be used to infer object depth. Since there are fewer keypoints than pixels, the consequent disparity or depth map is “sparse”. Keypoints may comprise any data structure which can be consistently replicated from a portion of an image and thereby permit unique identification of the image portion. The keypoint may be associated with a position in the image. The keypoint’s unique determination permits the keypoints to be identified from similar or identical portions in a second image. In some embodiments, keypoints may comprise Scale Invariant Feature Transform (SIFT) keypoints, or keypoints of a similar feature generation module. In some embodiments, the system may reuse machine vision components preexisting in the general processor 113 or subprocessors to determine keypoints. For example, high pass filtering blocks may be reused for keypoint detection. Alternatively, software libraries for performing machine vision operations stored in memory 103 may be used to generate keypoints. In this manner, certain implementations may economically take advantage of functionality associated with

other applications to generate keypoints for performing autofocus. Alternative means for determining a plurality of keypoints, such as feature generation modules employing algorithms other than SIFT, are described in greater detail below.

[0042] Figure 7 is a flow diagram depicting an autofocus process 700 for stereoscopic image capture, which may be implemented by certain of the present embodiments. The process begins 701 by acquiring, or receiving, at least a pair of stereoscopic images 702. Certain embodiments contemplate cropping a region of interest from each of the images, so as to reduce computation time. Once the images are received, the system determines keypoints from the first image 703. As mentioned, in some embodiments these keypoints may be determined using SIFT or other feature detection hardware, firmware, or software. In some embodiments, the system may also determine keypoints in the second image. The system may then correlate keypoints from the first image with pixel regions (such as a particular pixel position) in the second image 704. A “keypoint correlation” software, firmware, or hardware module may be configured to perform this operation. Certain portions of the operation may distributed among other modules (firmware, hardware, or software), creating other means for correlating keypoints. This operation may serve to identify the same image region in the first image in the second image.

[0043] Disparities D may then be calculated between each keypoint position of the first image and the correlated pixel positions of the second image 705. Where keypoints have been calculated for both images, the disparities between the keypoints may be determined by subtracting the relative positions of each of the correlated keypoints. The disparities may then be organized as a disparity histogram similar to 601 of Figure 6 and a corresponding depth histogram, similar to histogram 602, may be determined. The depth histogram may then be used to determine the optimal focal depth for the sensor based on the selected statistic. A “disparity determination” software, firmware, or hardware module may be configured to perform this operation. Certain portions of the operation may distributed among other modules (firmware, hardware, or software), creating other means for determining disparities.

[0044] In the embodiments implementing process 700, to improve computation efficiency, the process 700 determines the statistic (in this case, the average) of the disparities 706 rather than converting to a depth for each disparity and then determining the average of the depths. Only the depth of the single statistical value

need then be determined 707 with reference to a graph similar to that of Figure 4. This depth may then be used as the new camera focus depth 708 during subsequent image captures. As mentioned, other embodiments may instead convert each of the disparities to a depth and then average the depths. Other embodiments may alternatively take the mean, median, or some other statistic to determine the desired focal depth. A “depth determination” software, firmware, or hardware module may be configured to perform this operation. The module may operate in conjunction with a disparity histogram. Certain portions of the operation may distributed among other modules (firmware, hardware, or software), creating other means for determining a focal depth. Once the focus depth has been determined, image sensor 101a may be adjusted. Sensor 101b may also be adjusted by the processor 113, or sensor 101b may track the focal depth of sensor 101a independently. As mentioned above, in certain embodiments, only a single sensor may be adjusted based upon the determined focal depth.

[0045] In variations of these embodiments, the system may alternatively use information from the scene to determine the depth, rather than simply take the average of the disparities. For example, in lieu of taking the average, keypoint disparities may be weighted based upon their presence in a single object and upon lighting conditions. For example, the histograms may be enhanced by weighting each point from the histogram by the focal quality associated with a certain focus distance. In certain camera arrangements if the focal point was set to 3 meters, objects between 2m and infinity may have good focus, objects between 1m-2m may have fair focus, and objects between 0.5-1m may have bad focus. The histograms of Figure 6 would accordingly be weighted so that a preferred focal range was selected more often than the other ranges. This may be referred to as a “region weighted saliency” in certain embodiments. In other variations, frequency information from the image may be incorporated into the keypoints selection. Objects comprising textures may generate more keypoints than objects without textures or with little texture, and thereby affect the average. Accordingly, keypoints associated with textured objects may receive different weights from non-textured objects. In one variation, regions within a texture may be detected and these regions then used so as to lower the weight for keypoints in that region.

[0046] When capturing a stereoscopic movie, the process 700 may be applied to a single frame, i.e. a single pair of stereoscopic images. The determined focal depth may then be used by the image sensors during subsequent image captures until the

camera arrangement or scene is modified so as to necessitate reassessment of the proper focus. The operation 700 therefore has the benefit that there need not be any state dependency. That is, a traditional auto-focus system would need to periodically time-out the movie capture process and capture multiple focal depths to reassess the focus. Process 700, in contrast, may be “instantaneous” in that it produces no frame delay. This facilitates seamless focus tracking. The process may further guarantee system stability, as there is no feedback (or dependency) between the current focus position and the focus position estimation. Additionally, since the focus operation may be accomplished with a single frame, the user’s hand motion will be less likely to generate any blur.

[0047] As mentioned, the keypoints generated in steps 703 and 704 of process 700 may comprise any data structure which can assign an identity to a portion of the image and be consistently recognized when applied to the second image. As mentioned, in some embodiments the keypoints may comprise Scale Invariant Feature Transform (SIFT) keypoints generated from a SIFT feature generation module. Figure 8 depicts another possible process for generating keypoints in a feature generation module.

[0048] The process 800 begins 801 by receiving 802 one of the pair of raw stereoscopic images. The image may then be subsampled 803, possibly to improve the algorithm’s robustness to noise and to decrease computational demands. The image may then be passed through a horizontal high pass filter 804. In some embodiments, the filter may comprise a 3x4 kernel with a response given by

$$h = \begin{bmatrix} -1 & -1 & 1 & 1 \\ -1 & -1 & 1 & 1 \\ -1 & -1 & 1 & 1 \end{bmatrix}$$

[0049] The process may then calculate the power 805 of the image, i.e. by taking the square of each value. Finally, the process may threshold the values 806 to eliminate noise and low-power values. From among the remaining values which exceeded the threshold, the system will identify “maximum values”. In some embodiments the maximum values may be those image portions which exceeded the threshold, while in other embodiments the maximum values may be defined relative to their local neighbors. For example, the delta between neighboring pixels exceeding the threshold may be used to identify a maximum value. The identified maximum values

represent keypoint positions which may be used for the disparity determination steps described above. The system may store these keypoint positions 807 before ending 808. By subsampling and thresholding the image in this manner the computation time required to determine pixel locations which may serve as keypoints may be reduced. As a keypoint in these embodiments comprises a pixel position, the pixel position may occasionally be referred to as a “keypoint”. However, one will readily recognize variations, wherein keypoints comprise both a position and an array of the neighboring pixel values and positions. Keypoints may also refer to frequency content of an image portion or to gradients in pixel value, rather than to pixel values or pixel positions directly. SIFT keypoints, for example, may comprise a vector indicating the pixel gradient.

[0050] Once the keypoints in the first image have been determined, it may still remain to correlate the keypoints with positions in the second image so that the disparities between image portions may be determined. A possible method for determining whether keypoints, such as the keypoints generated by the process 800 of Figure 8, are correlated with a position in an image will now be described with respect to Figure 9. Figure 9 depicts stereoscopic image pair 900a and 900b. Based on the configuration of the two sensors used to capture images 900a and 900b, a search region 902 around each keypoint 901 may be determined. The search region 902 may specify the maximum distance a keypoint may be displaced in the left and right images (i.e., a maximum expected disparity) as a consequence of the capture device configuration. In some embodiments, as the image sensors may lack vertical disparity, keypoints 901 may generally lie on vertical edges in the scene.

[0051] Search region 902 is located in the same absolute position of each of images 900a and 900b. In Figure 9, search region 902 comprises a rectangle having the height of a single pixel, since it is assumed that no vertical disparity between images 900a and 900b exists (i.e. only pixels in a same row are considered). The height of the region 902 may be increased in relation to the amount of vertical disparity that may be present between images 900a and 900b. The system may iterate through each pixel in the search region 902 of the second image 900b and determine the pixel’s correspondence with the portions of image 900a surrounding the keypoint 901. This may be accomplished in some embodiments using a correlation metric described in further detail below.

[0052] In certain embodiments, keypoints may have been determined for both image 900a and image 900b, rather than simply for image 900a. When attempting to correlate keypoints between the images, the system may identify keypoints within the search region 902 in the second image 900b. If only one keypoint of image 900b is found in search region 902, this keypoint may be correlated with the keypoint 901 from the first image. Where more than one keypoint of image 900b is present in the region 902, the system may apply a correlation metric to each keypoint of image 900b in the search region 902 to determine which keypoint best corresponds with the keypoint 901 from the first image. As with metrics applied when keypoints are taken for only one image, the metric may consider pixel values neighboring the pixel positions of keypoints 901 and 901b to verify that keypoints 901 and 901b are more likely to refer to the same portion of the scene in each of images 900a and 900b. Where keypoints are created for both images, it may be necessary only to iterate between keypoints in the region 902 rather than between each pixel within the region.

[0053] In the embodiments described above, the system iterates through certain of the pixels in the search region 902 corresponding to the determined keypoint 901. The system may apply a correlation metric to each pixel in the region 902. The pixel in the region 902 having the maximum correlation with the region surrounding the position of keypoint 901 may then be correlated with the keypoint 901. The computational cost to iterate through each pixel of image 900b in the range 902 may be less than the cost to compute keypoints for all of image 900b and to then determine the correlations between each keypoint. In some embodiments, however, where only a few keypoints have been generated, the system may determine correlations between all the keypoints directly, rather than iterate between regions 902 associated with the keypoints of one image.

[0054] In certain embodiments, the correlation metric used to identify a keypoint or pixel position in image 900b corresponding to a keypoint in 900a may comprise the calculation of the mean square error for pixels surrounding a position in image 900b under consideration and the pixels surrounding the keypoint position in image 900a. That is, the mean square error of pixels neighboring keypoint 901 of the first image and the neighboring pixels for positions in search region 902 in image 900b may be used as a correlation metric. The mean square error may be calculated as:

$$R(\Delta) = \sum_{i=-3}^3 \sum_{j=-3}^3 (S_{left}(i+M, j+N) - S_{right}(i+M+\Delta, j+N))^2$$

where R is the mean squared error, S_{left} comprises the pixel values in the image 900a, S_{right} comprises the pixel values in the image 900b, M and N comprise the horizontal and vertical offsets into the image to the region 902 for the current keypoint or pixel position of region 902 under investigation, and Δ comprises the horizontal shift applied for the current position in the search region 902 (the first parameter to S_{left} and S_{right} is a column position/x-axis and the second a row position/y-axis). Although the mean squared error is within a 7x7 window in the above example, one may readily envision a range of window dimensions depending upon the resolution of the image and the subsampling applied. Furthermore, as sensors 101a and 101b are presumed not to have any vertical disparity in the above example, the search region 902 extends only horizontally and Δ appears only in the x-axis/column direction. More robust systems may compensate for errors in sensor positioning by increasing the height of the search region 902 and including a Δ in the vertical direction. As the image has been downsampled, sub pixel resolution may be determined using interpolation, such as polynomial interpolation, in some embodiments to facilitate more accurate determinations of the pixel in region 902 of image 900b correlated with the keypoint 901. That is, the displacement of sensor 101b relative to sensor 101a may not be an exact, integer number of pixels. Thus, particularly after sub-sampling, accurate correlation of the keypoint 901 may require including locations between pixels in the search region 902. The, position in region 902 of image 900b maximally correlated with keypoint 901 may fall between pixel positions, at an interpolated point.

[0055] The various illustrative logical blocks, modules, and circuits described in connection with the implementations disclosed herein may be implemented or performed with a general purpose processor, a digital signal processor (DSP), an application specific integrated circuit (ASIC), a field programmable gate array (FPGA) or other programmable logic device, discrete gate or transistor logic, discrete hardware components, or any combination thereof designed to perform the functions described herein. A general purpose processor may be a microprocessor, but in the alternative, the processor may be any conventional processor, controller, microcontroller, or state machine. A processor may also be implemented as a combination of computing devices, e.g., a combination of a DSP and a microprocessor, a plurality of

microprocessors, one or more microprocessors in conjunction with a DSP core, or any other such configuration.

[0056] The steps of a method or process described in connection with the implementations disclosed herein may be embodied directly in hardware, in a software module executed by a processor, or in a combination of the two. A software module may reside in RAM memory, flash memory, ROM memory, EPROM memory, EEPROM memory, registers, hard disk, a removable disk, a CD-ROM, or any other form of non-transitory storage medium known in the art. An exemplary computer-readable storage medium is coupled to the processor such the processor can read information from, and write information to, the computer-readable storage medium. In the alternative, the storage medium may be integral to the processor. The processor and the storage medium may reside in an ASIC. The ASIC may reside in a user terminal, camera, or other device. In the alternative, the processor and the storage medium may reside as discrete components in a user terminal, camera, or other device.

[0057] Headings are included herein for reference and to aid in locating various sections. These headings are not intended to limit the scope of the concepts described with respect thereto. Such concepts may have applicability throughout the entire specification.

[0058] The previous description of the disclosed implementations is provided to enable any person skilled in the art to make or use the present invention. Various modifications to these implementations will be readily apparent to those skilled in the art, and the generic principles defined herein may be applied to other implementations without departing from the spirit or scope of the invention. Thus, the present invention is not intended to be limited to the implementations shown herein but is to be accorded the widest scope consistent with the principles and novel features disclosed herein.

We claim:

1. A method in an electronic device for determining a focal depth for an image sensor, comprising:
 - receiving a first image associated with a first viewpoint;
 - receiving a second image associated with a second viewpoint;
 - determining a first plurality of keypoints based on the first image;
 - correlating keypoints from the first plurality of keypoints with positions in the second image;
 - determining a plurality of disparities associated with each of the first plurality of keypoints; and
 - determining a focal depth based upon the plurality of disparities, the position of the first viewpoint and the position of the second viewpoint.
2. The method of Claim 1, further comprising determining a second plurality of keypoints based on the second image.
3. The method of Claim 2, wherein correlating keypoints from the first plurality of keypoints with positions in the second image comprises correlating keypoints from the first plurality of keypoints with keypoints from the second plurality of keypoints.
4. The method of Claim 1, wherein correlating keypoints from the first plurality of keypoints with positions in the second image comprises iterating over pixels within a search range in the second image.
5. The method of Claim 4, wherein correlating keypoints from the first plurality of keypoints with positions in the second image comprises determining the mean square error between pixels in the first image and the second image.
6. The method of Claim 1, wherein determining a first plurality of keypoints based on the first image comprises determining Scale Invariant Feature Transform (SIFT) keypoints based on the first image.
7. The method of Claim 1, wherein determining a first plurality of keypoints based on the first image comprises sub-sampling the first image, applying a

high-pass filter to the first image, calculating the power of the first image, and thresholding the first image.

8. The method of Claim 1, wherein correlating keypoints from the first plurality of keypoints with positions in the second image occurs in realtime.

9. The method of Claim 1, wherein the electronic device comprises a mobile phone.

10. A computer readable medium comprising instructions configured to cause a computer to perform the steps of:

receiving a first image associated with a first viewpoint;

receiving a second image associated with a second viewpoint;

determining a first plurality of keypoints based on the first image;

correlating keypoints from the first plurality of keypoints with positions in the second image;

determining a plurality of disparities associated with each of the first plurality of keypoints; and

determining a focal depth based upon the plurality of disparities, the position of the first viewpoint and the position of the second viewpoint.

11. The computer readable medium of Claim 10, further comprising determining a second plurality of keypoints based on the second image.

12. The computer readable medium of Claim 11, wherein correlating keypoints from the first plurality of keypoints with positions in the second image comprises correlating keypoints from the first plurality of keypoints with keypoints from the second plurality of keypoints.

13. The computer readable medium of Claim 10, wherein correlating keypoints from the first plurality of keypoints with positions in the second image comprises iterating over pixels within a search range in the second image.

14. The computer readable medium of Claim 13, wherein correlating keypoints from the first plurality of keypoints with positions in the second image comprises determining the mean square error between pixels in the first image and the second image.

15. The computer readable medium of Claim 10, wherein determining a first plurality of keypoints based on the first image comprises determining Scale Invariant Feature Transform (SIFT) keypoints based on the first image.

16. The computer readable medium of Claim 10, wherein determining a first plurality of keypoints based on the first image comprises sub-sampling the first image, applying a high-pass filter to the first image, calculating the power of the first image, and thresholding the first image.

17. The computer readable medium of Claim 10, wherein correlating keypoints from the first plurality of keypoints with positions in the second image occurs in realtime.

18. The computer readable medium of Claim 10, wherein the computer is located in a mobile phone.

19. A system for focusing a stereoscopic capture device, the system comprising:

- a first image sensor configured to generate a first image associated with a first viewpoint;

- a second image sensor configured to generate a second image associated with a second viewpoint;

- a feature generation module configured to determine a first plurality of keypoints based on the first image;

- a keypoint correlation module configured to correlate keypoints from the first plurality of keypoints with positions in the second image;

- a disparity determination module configured to determine a plurality of disparities associated with each of the first plurality of keypoints; and

- a depth determination module configured to determine a focal depth based upon the plurality of disparities, the position of the first viewpoint and the position of the second viewpoint.

20. The system of Claim 19, wherein the feature generation module is configured to determine a second plurality of keypoints based on the second image.

21. The system of Claim 20, wherein the module configured to correlate keypoints is configured to correlate keypoints from the first plurality of keypoints with keypoints from the second plurality of keypoints.

22. The system of Claim 19, wherein the module configured to correlate keypoints is configured to iterate over pixels within a search range in the second image.

23. The system of Claim 22, wherein correlating keypoints from the first plurality of keypoints with positions in the second image comprises determining the mean square error between pixels in the first image and the second image.

24. The system of Claim 19, wherein the feature generation module is configured to determine Scale Invariant Feature Transform (SIFT) keypoints based on the first image.

25. The system of Claim 19, wherein the feature generation module is configured to sub-sample the first image, apply a high-pass filter to the first image, calculate the power of the first image, and threshold the first image.

26. The system of Claim 19, wherein the module configured to correlate keypoints correlates keypoints from the first plurality of keypoints with positions in the second image in realtime.

27. The system of Claim 19, wherein the stereoscopic capture device is located on a mobile phone.

28. The system of Claim 19, wherein the module configured to determine a focal depth comprises a disparity histogram.

29. A system for focusing a stereoscopic capture device, the system comprising:

means for receiving a first image associated with a first viewpoint;

means for receiving a second image associated with a second viewpoint;

means for determining a first plurality of keypoints based on the first image;

means for correlating keypoints from the first plurality of keypoints with positions in the second image;

means for determining a plurality of disparities associated with each of the first plurality of keypoints; and

means for determining a focal depth based upon the plurality of disparities, the position of the first viewpoint and the position of the second viewpoint.

30. The system of Claim 29, wherein the means for receiving a first image comprises a first sensor, the means for receiving a second image comprises a second sensor, the means for determining a first plurality of keypoints comprises a feature generation module, the means for correlating comprises a keypoint correlation module,

the means for determining a plurality of disparities comprises a disparity determination module, and the means for determining a focal depth comprises a depth determination module.

31. The system of Claim 29, wherein the means for determining a first plurality of keypoints is configured to determine a second plurality of keypoints based on the second image.

32. The system of Claim 31, wherein the means for correlating keypoints from the first plurality of keypoints with positions in the second image is configured to correlate keypoints from the first plurality of keypoints with keypoints from the second plurality of keypoints.

33. The system of Claim 29, wherein the means for correlating keypoints from the first plurality of keypoints with positions in the second image is configured to iterate over pixels within a search range in the second image.

34. The system of Claim 33, wherein the means for correlating keypoints from the first plurality of keypoints with positions in the second image is configured to determine the mean square error between pixels in the first image and the second image.

35. The system of Claim 29, wherein the means for determining a first plurality of keypoints is configured to determine Scale Invariant Feature Transform (SIFT) keypoints based on the first image.

36. The system of Claim 29, wherein the means for determining a first plurality of keypoints is configured to sub-sample the first image, apply a high-pass filter to the first image, calculate the power of the first image, and threshold the first image.

37. The system of Claim 29, wherein the means for correlating keypoints correlates the keypoints from the first plurality of keypoints with positions in the second image in realtime.

38. The system of Claim 29, wherein the stereoscopic capture device is located on a mobile phone.

1/9

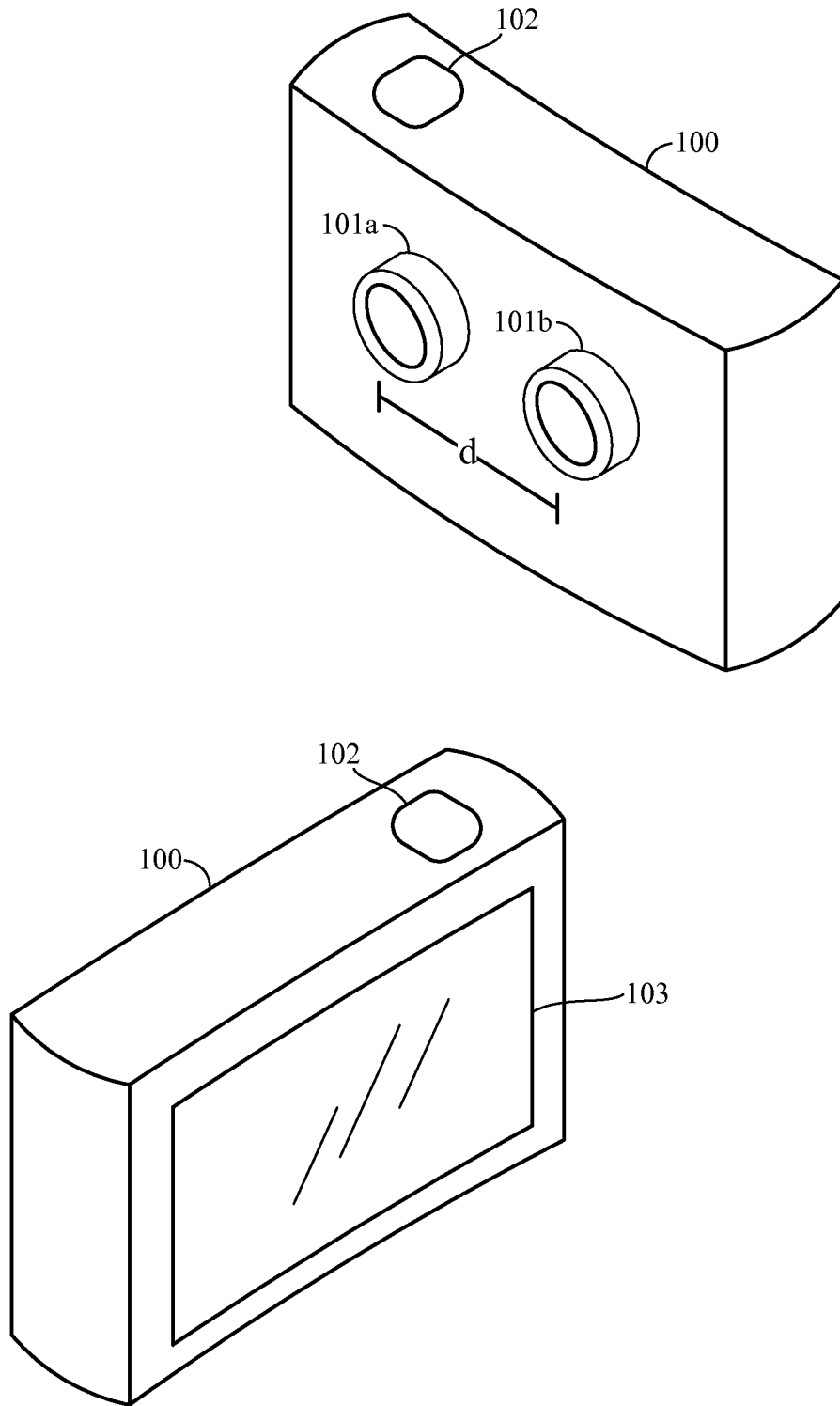


FIG. 1

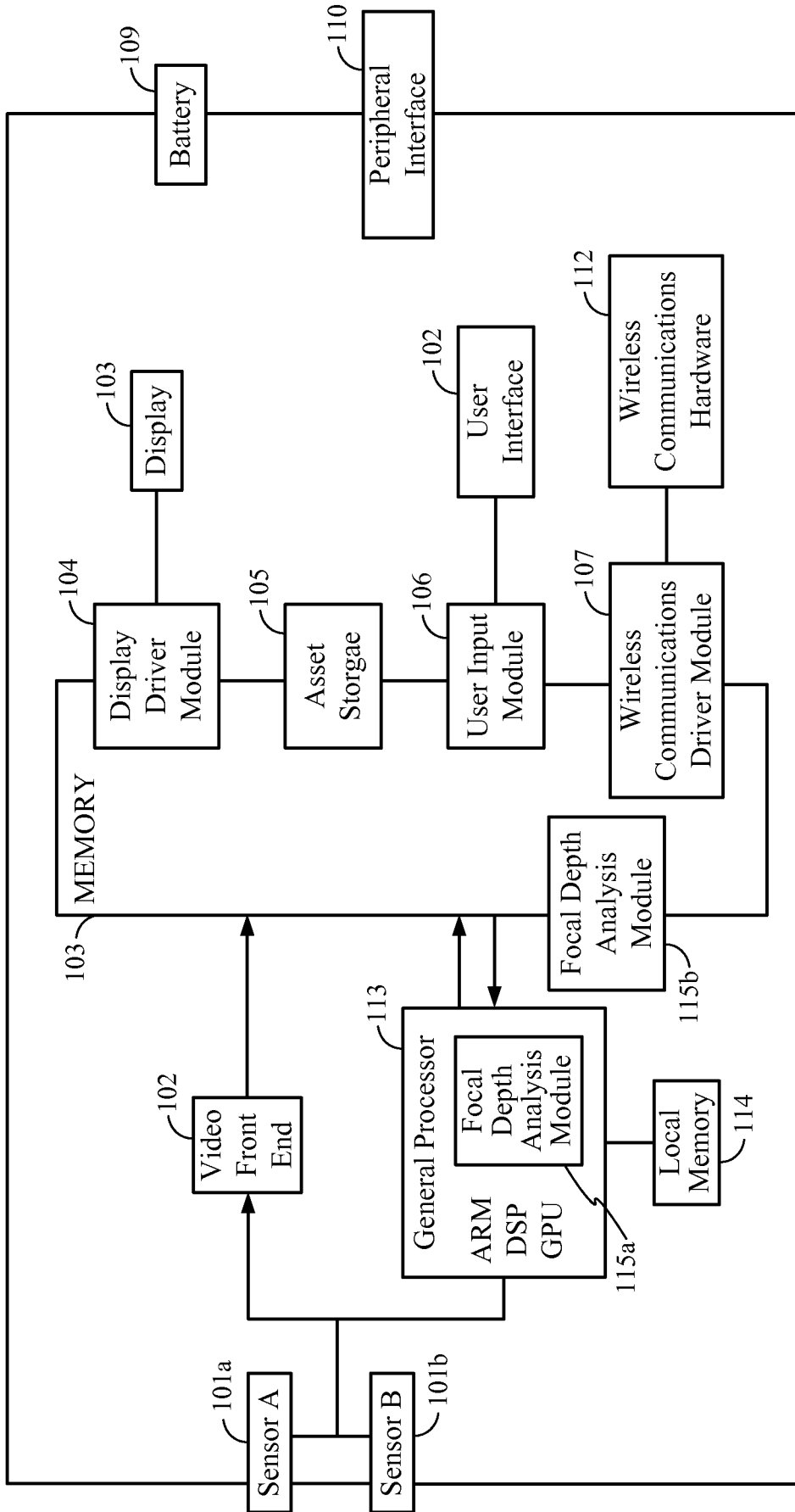


FIG. 2

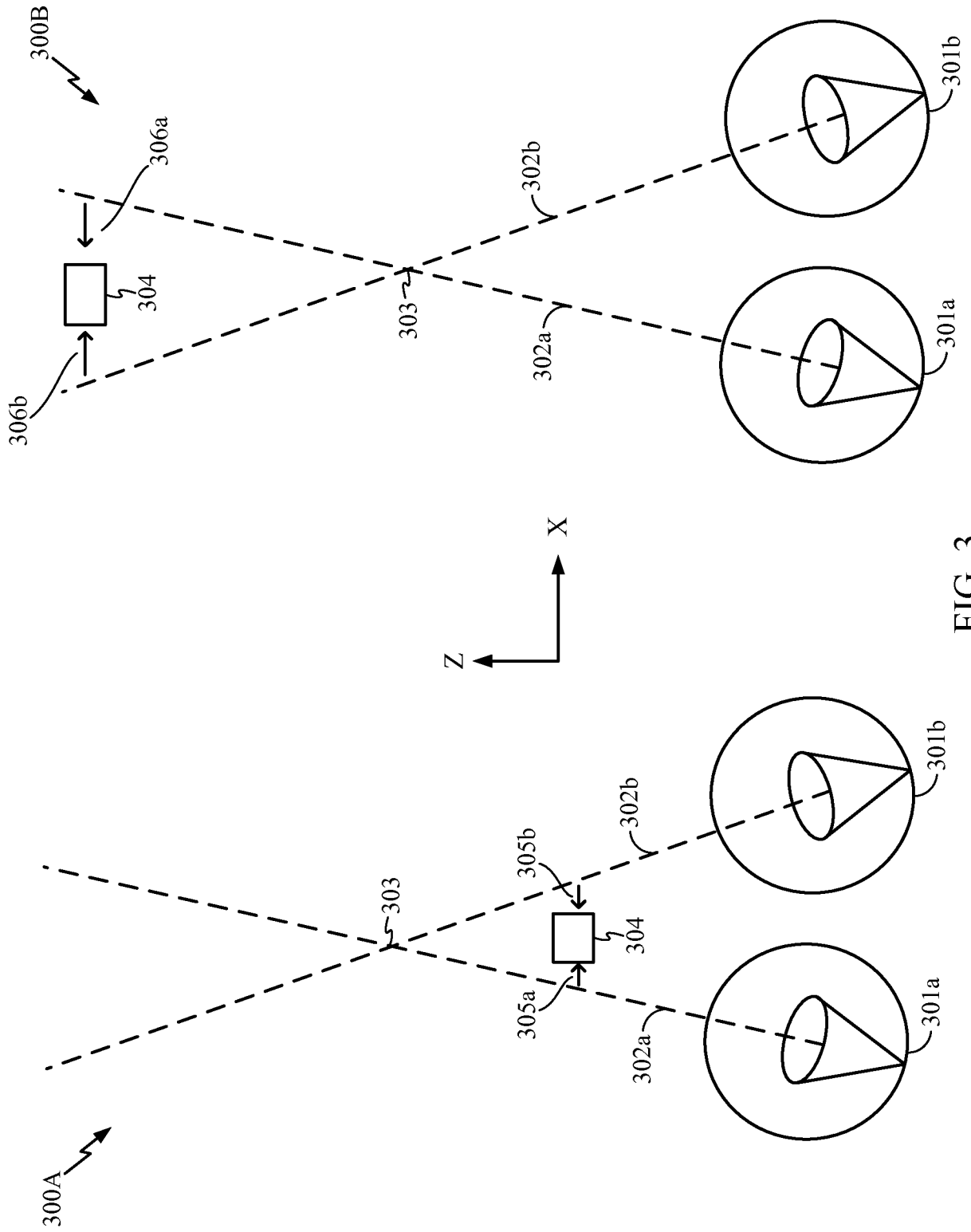


FIG. 3

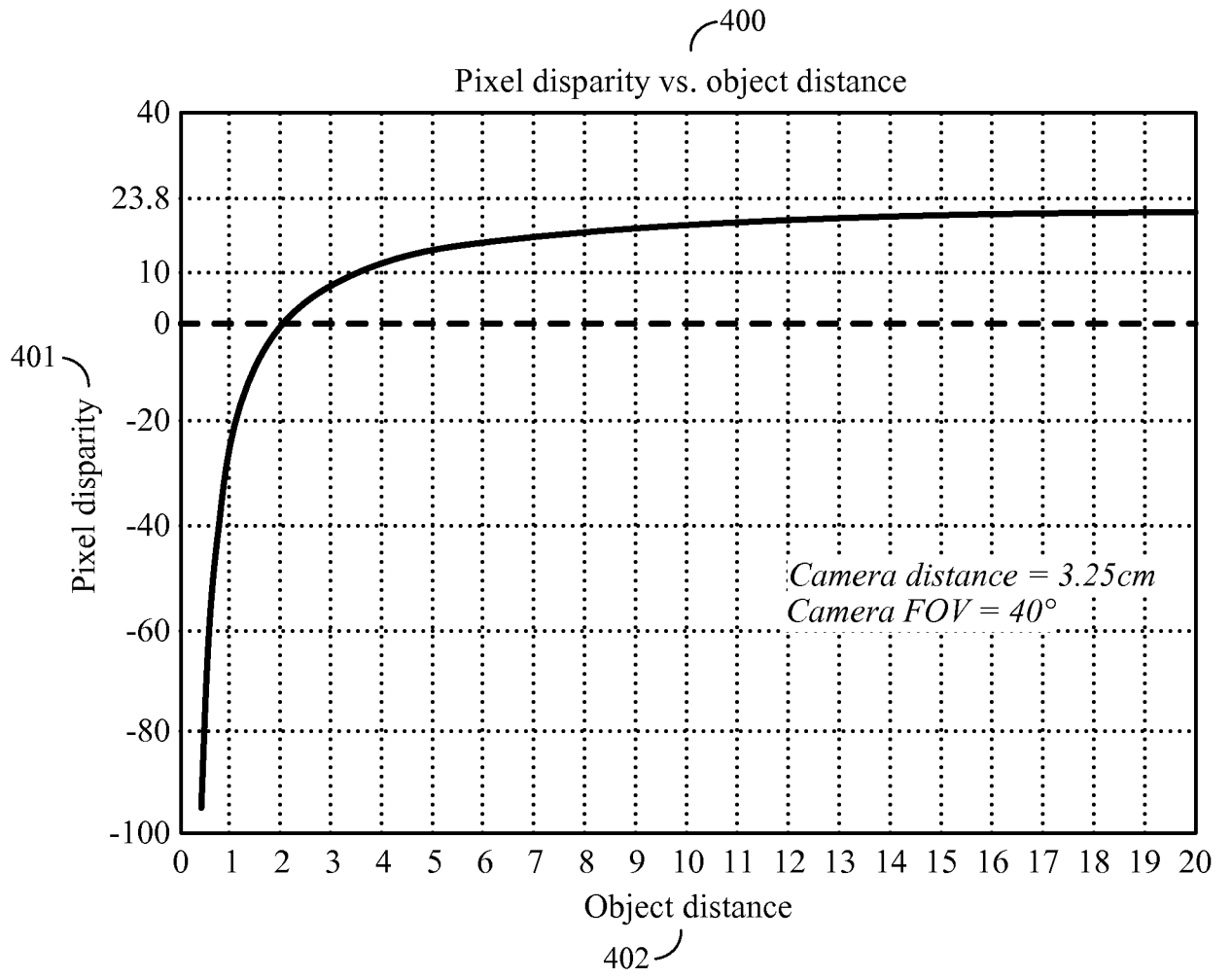


FIG. 4

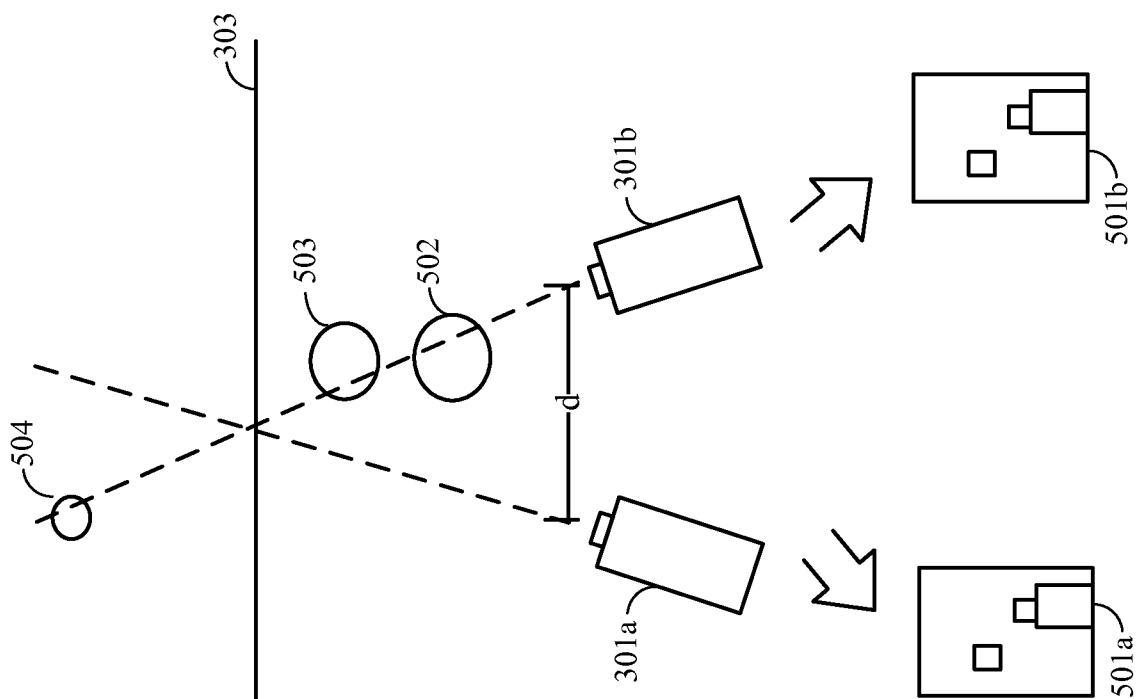


FIG. 5A

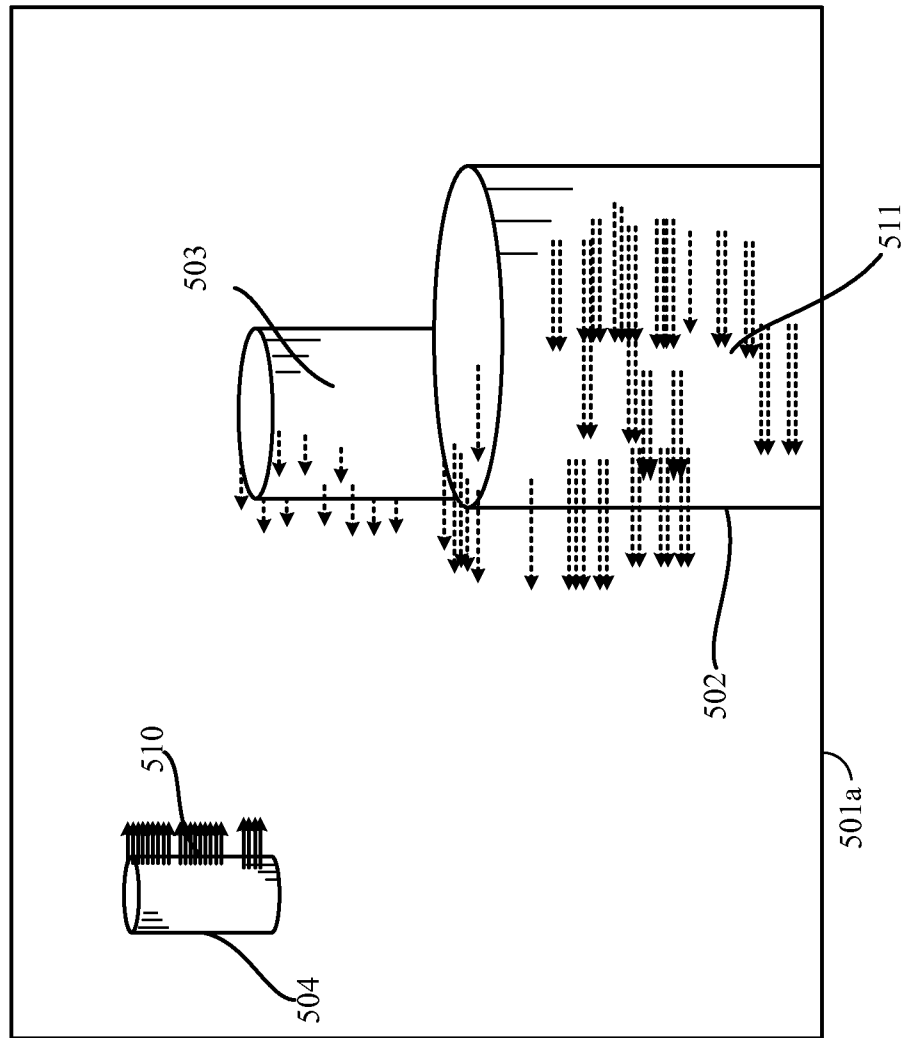


FIG. 5B

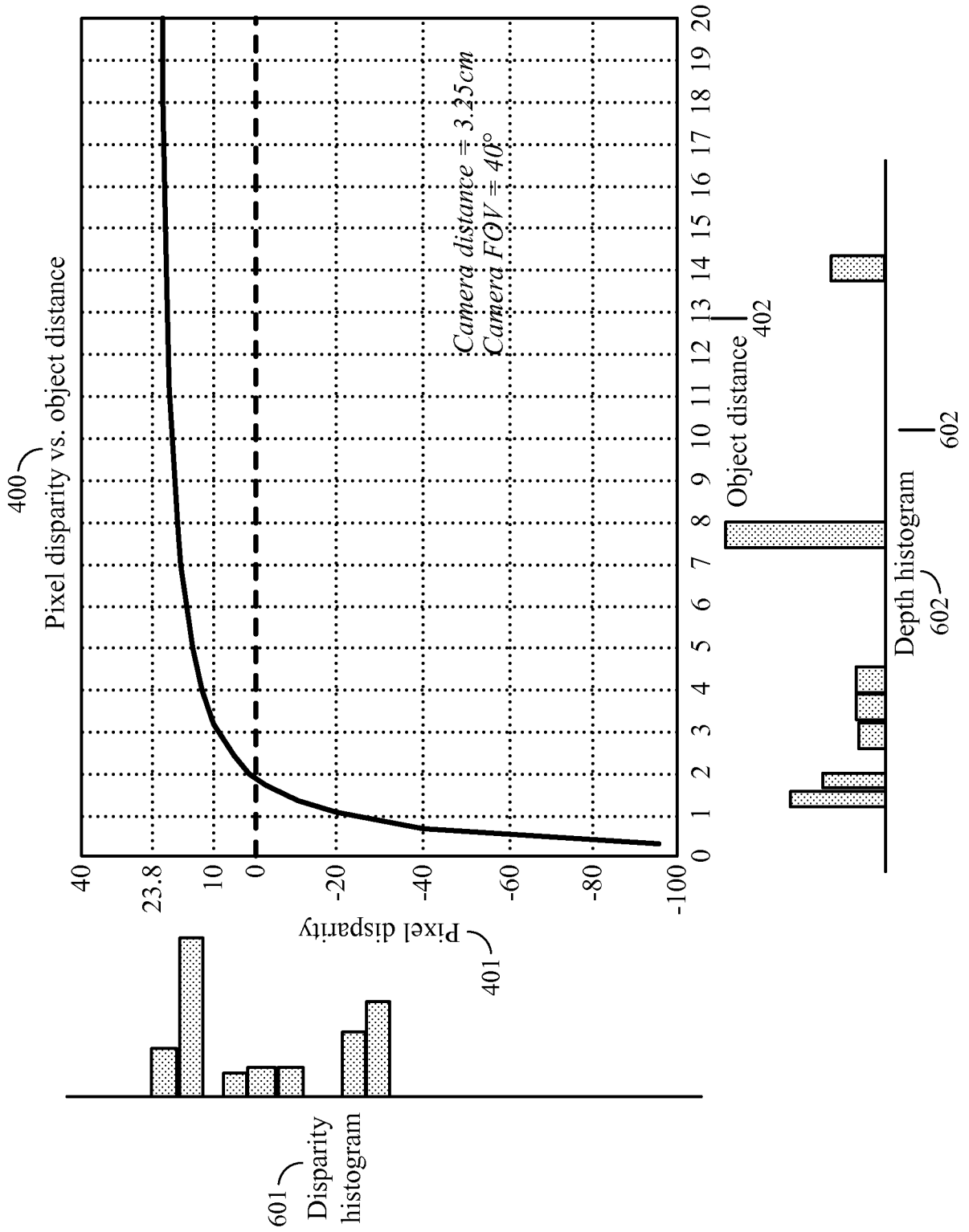


FIG. 6

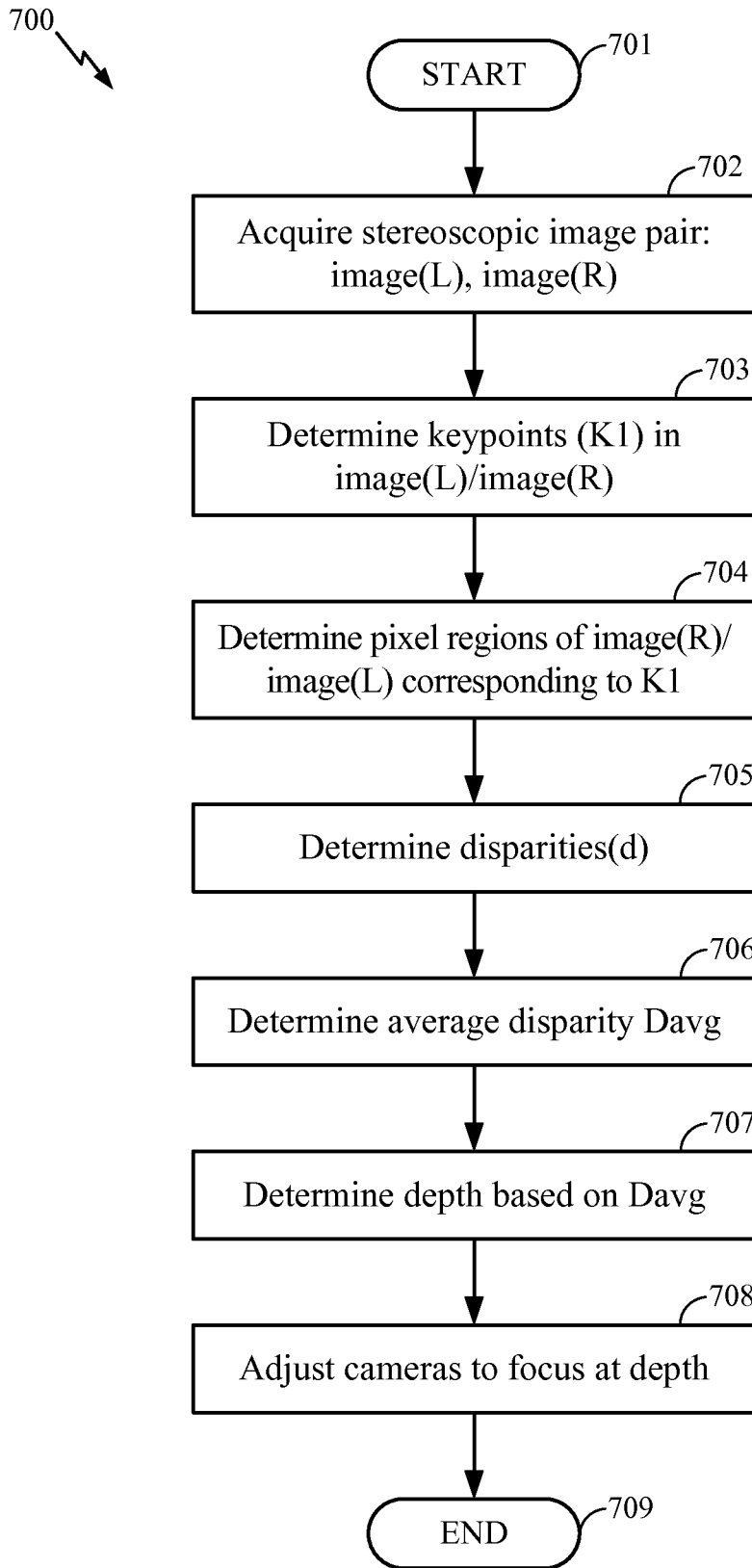


FIG. 7

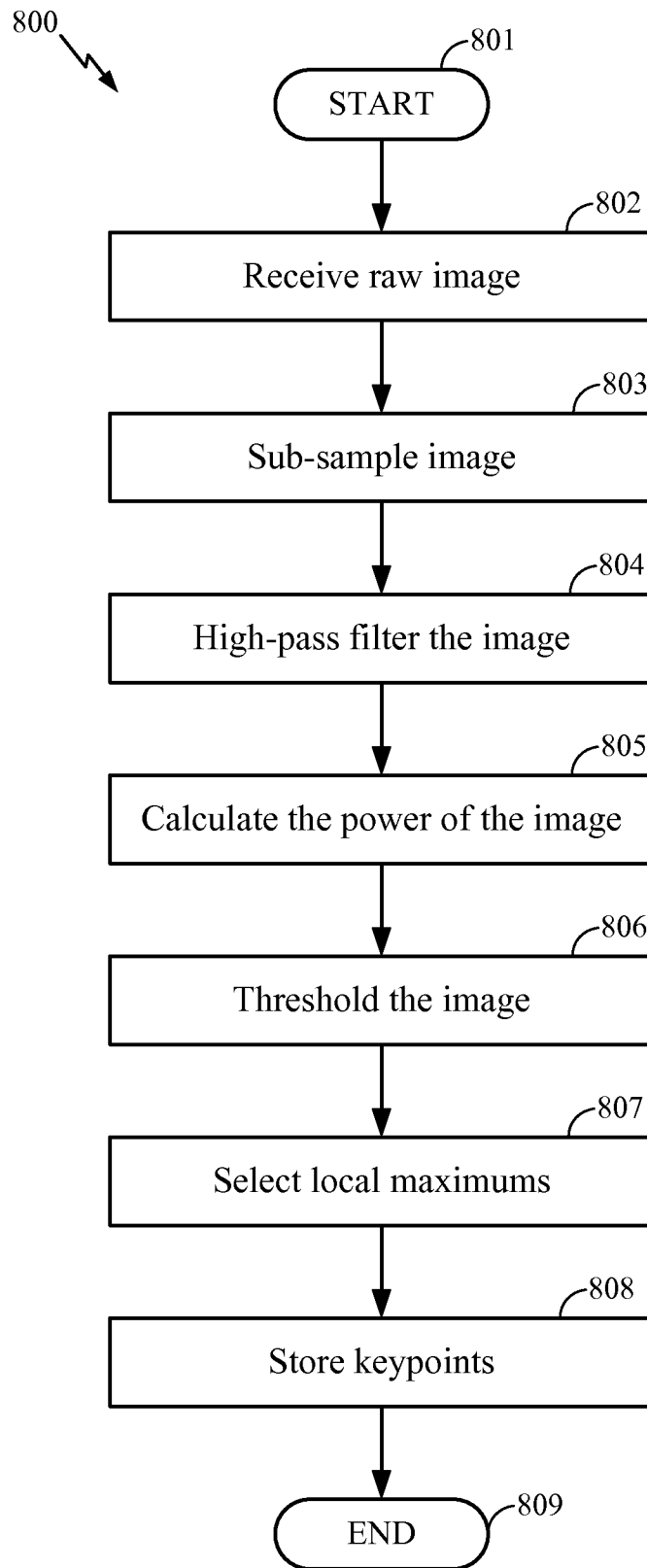


FIG. 8

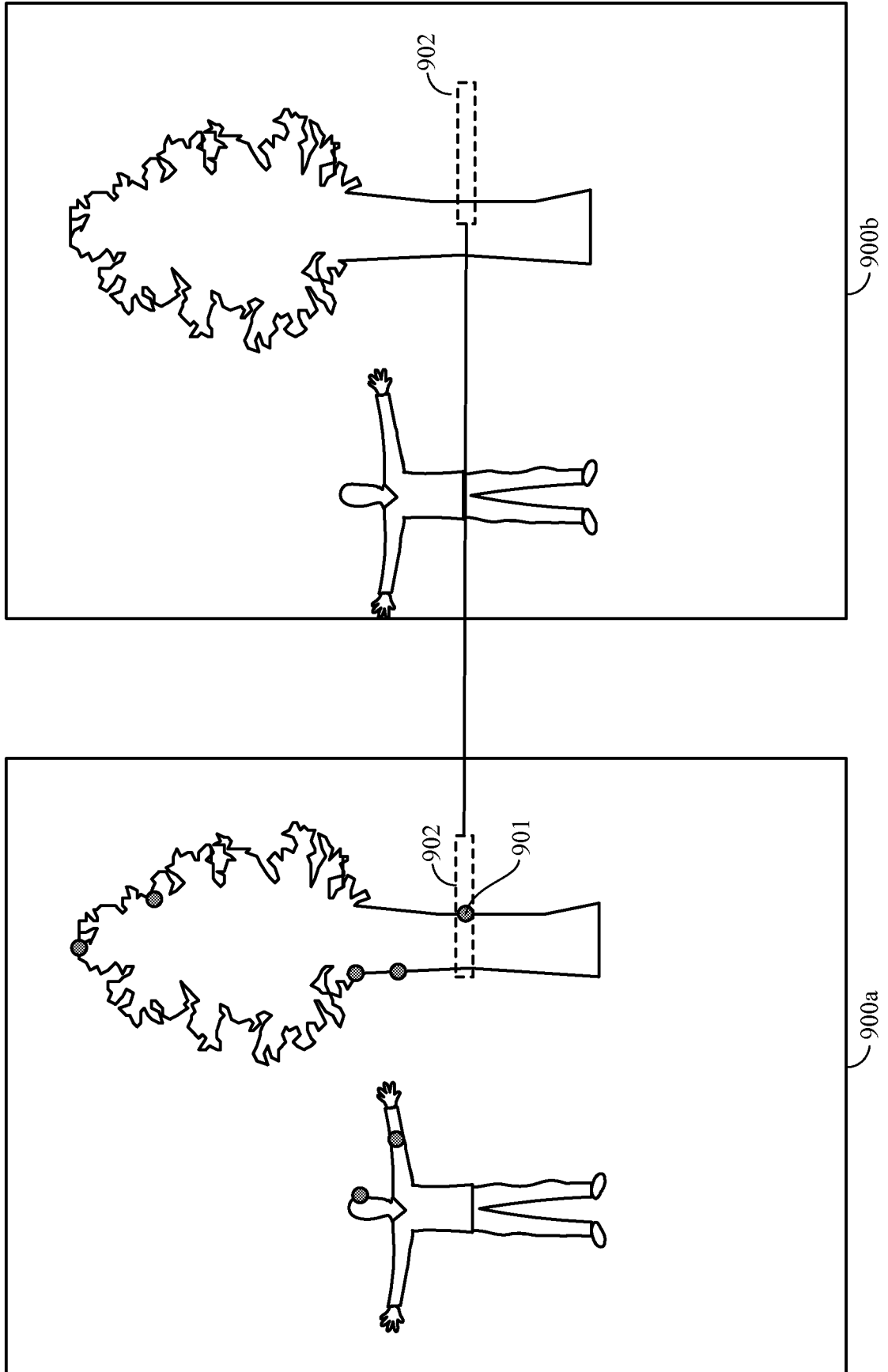


FIG. 9

INTERNATIONAL SEARCH REPORT

International application No
PCT/US2011/047126

A. CLASSIFICATION OF SUBJECT MATTER

INV. H04N5/232 H04N13/02 H04N13/00
ADD.

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)
H04N

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

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

EPO-Internal

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	JP 8 194274 A (OLYMPUS OPTICAL CO) 30 July 1996 (1996-07-30) abstract; figures 1,20,22 -----	1-38
A	US 2008/259172 A1 (TAMARU MASAYA [JP]) 23 October 2008 (2008-10-23) abstract; figure 2 -----	1-38
A	US 2008/218612 A1 (BORDER JOHN N [US] ET AL) 11 September 2008 (2008-09-11) paragraph [0108]; figures 3,8 -----	1-38
A	US 2006/215903 A1 (NISHIYAMA MANABU [JP]) 28 September 2006 (2006-09-28) paragraphs [0011] - [0012]; figure 1 -----	1-38
A	JP 2007 147457 A (TOPCON CORP) 14 June 2007 (2007-06-14) abstract; figure 4 -----	1-38

Further documents are listed in the continuation of Box C.

See patent family annex.

* Special categories of cited documents :

- "A" document defining the general state of the art which is not considered to be of particular relevance
- "E" earlier document but published on or after the international filing date
- "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)
- "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

- "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
- "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
- "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.
- "&" document member of the same patent family

Date of the actual completion of the international search

17 October 2011

Date of mailing of the international search report

26/10/2011

Name and mailing address of the ISA/

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040,
Fax: (+31-70) 340-3016

Authorized officer

Zakharian, Andre

INTERNATIONAL SEARCH REPORT

Information on patent family members

International application No

PCT/US2011/047126

Patent document cited in search report	A	Publication date	Patent family member(s)	Publication date
JP 8194274	A	30-07-1996	NONE	
US 2008259172	A1	23-10-2008	JP 4582423 B2	17-11-2010
			JP 2008271241 A	06-11-2008
US 2008218612	A1	11-09-2008	CN 101637019 A	27-01-2010
			EP 2119222 A1	18-11-2009
			JP 2010524279 A	15-07-2010
			WO 2008112054 A1	18-09-2008
US 2006215903	A1	28-09-2006	JP 4177826 B2	05-11-2008
			JP 2006268345 A	05-10-2006
JP 2007147457	A	14-06-2007	NONE	