



US 20160189414A1

(19) **United States**

(12) **Patent Application Publication**
BAKER et al.

(10) **Pub. No.: US 2016/0189414 A1**

(43) **Pub. Date: Jun. 30, 2016**

(54) **AUTOCAPTIONING OF IMAGES**

Publication Classification

(71) Applicant: **Microsoft Technology Licensing, LLC**,
Redmond, WA (US)

(51) **Int. Cl.**
G06T 11/60 (2006.01)
G06F 17/24 (2006.01)

(72) Inventors: **Simon BAKER**, Palo Alto, CA (US);
Krishnan RAMNATH, Redmond, WA (US)

(52) **U.S. Cl.**
CPC **G06T 11/60** (2013.01); **G06F 17/24**
(2013.01)

(73) Assignee: **Microsoft Technology Licensing, LLC**,
Redmond, WA (US)

(57) **ABSTRACT**

(21) Appl. No.: **15/063,323**

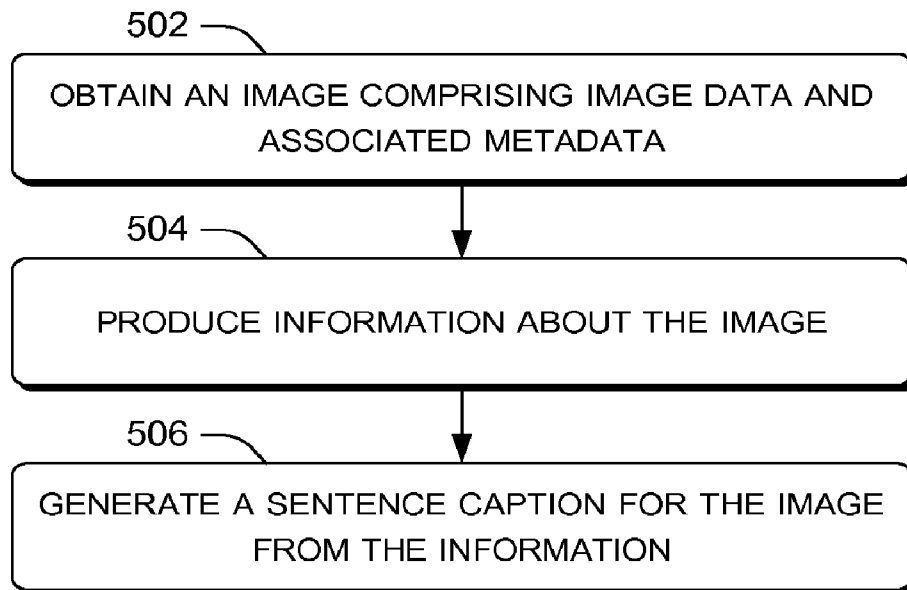
(22) Filed: **Mar. 7, 2016**

The description relates to sentence autocaptioning of images. One example can include a set of information modules and a set of sentence generation modules. The set of information modules can include individual information modules configured to operate on an image or metadata associated with the image to produce image information. The set of sentence generation modules can include individual sentence generation modules configured to operate on the image information to produce a sentence caption for the image.

Related U.S. Application Data

(63) Continuation of application No. 13/654,419, filed on Oct. 18, 2012, now Pat. No. 9,317,531.

METHOD 500



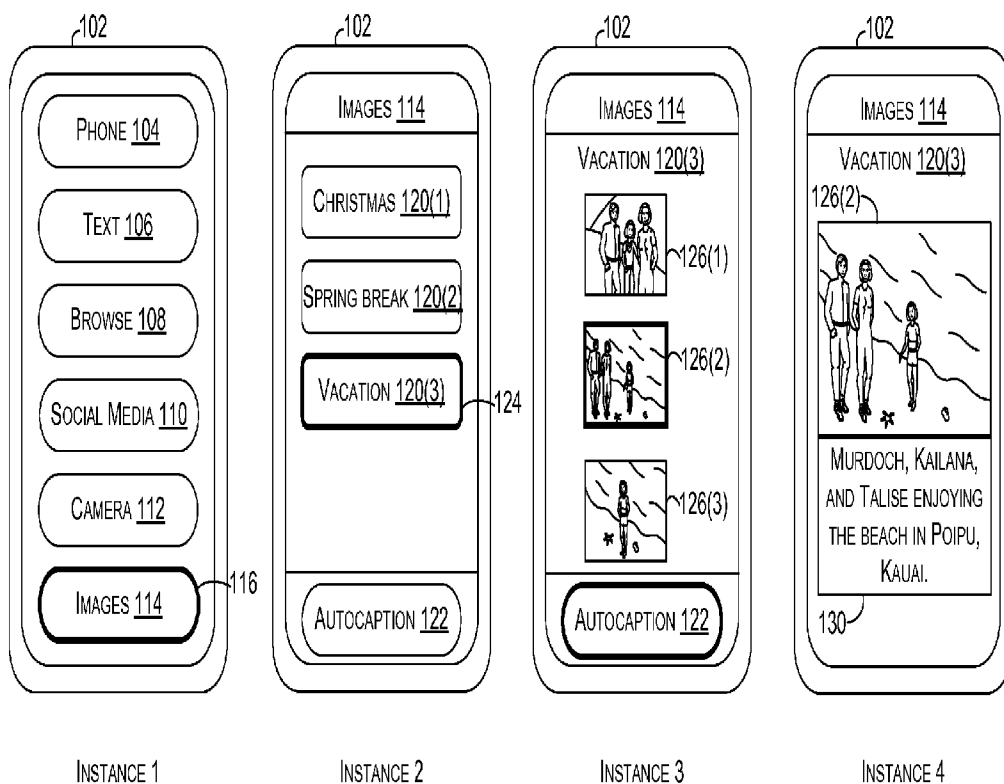


FIG. 1A

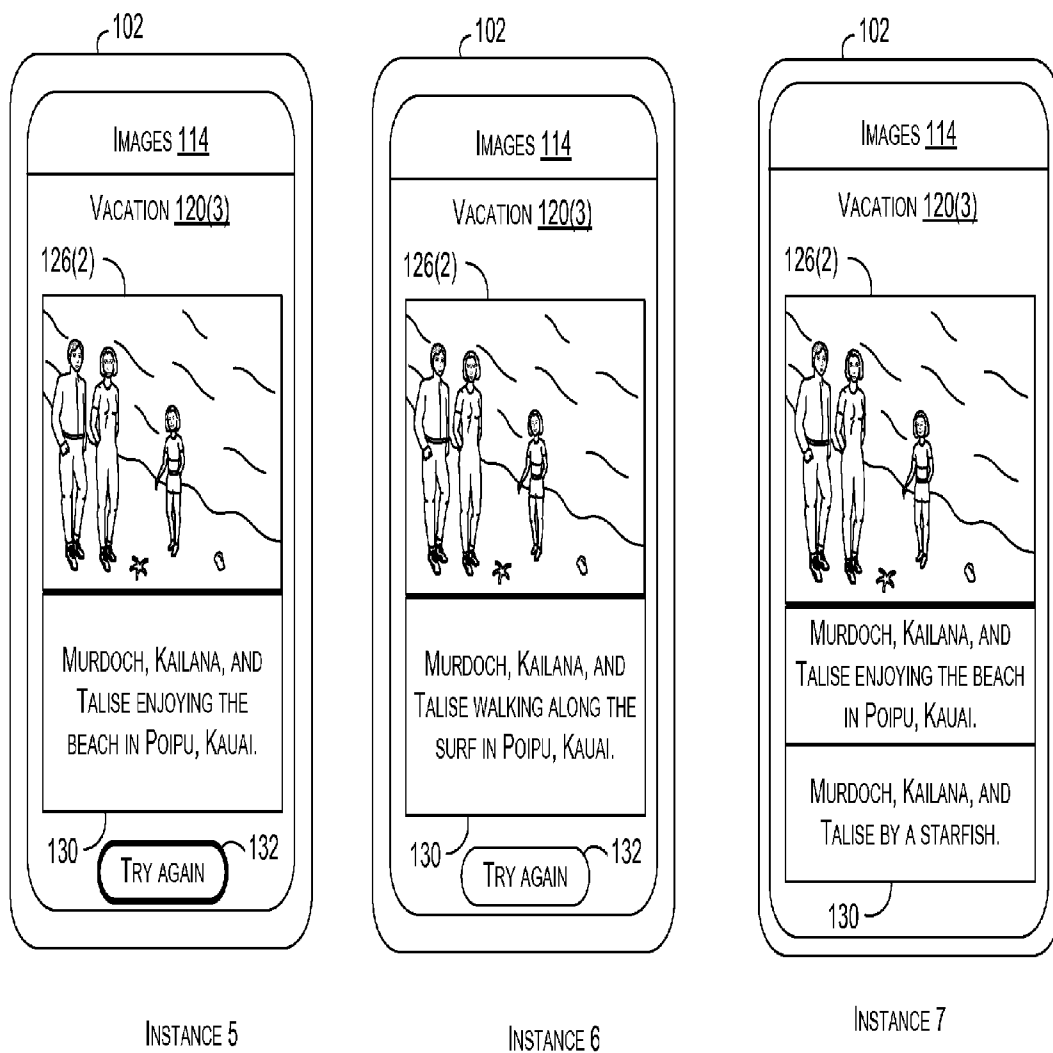


FIG. 1B

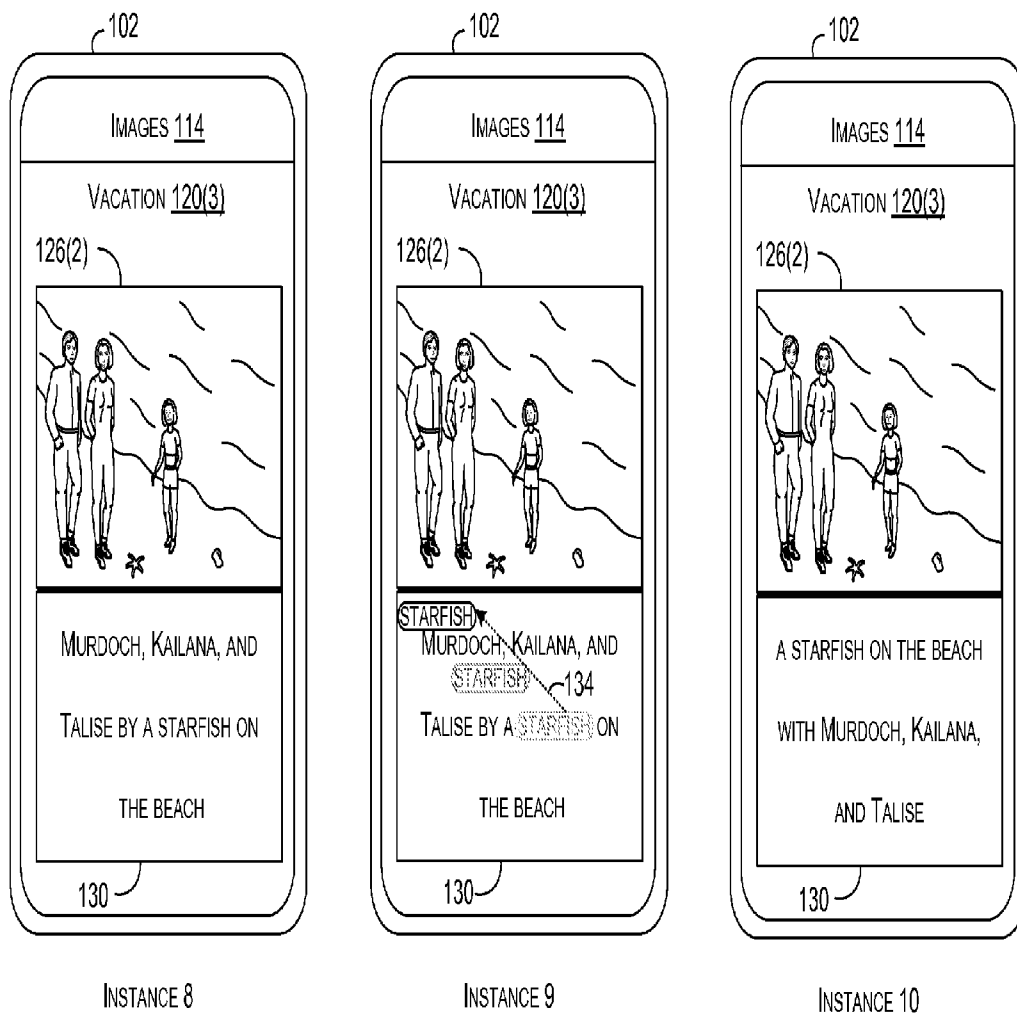


FIG. 1C

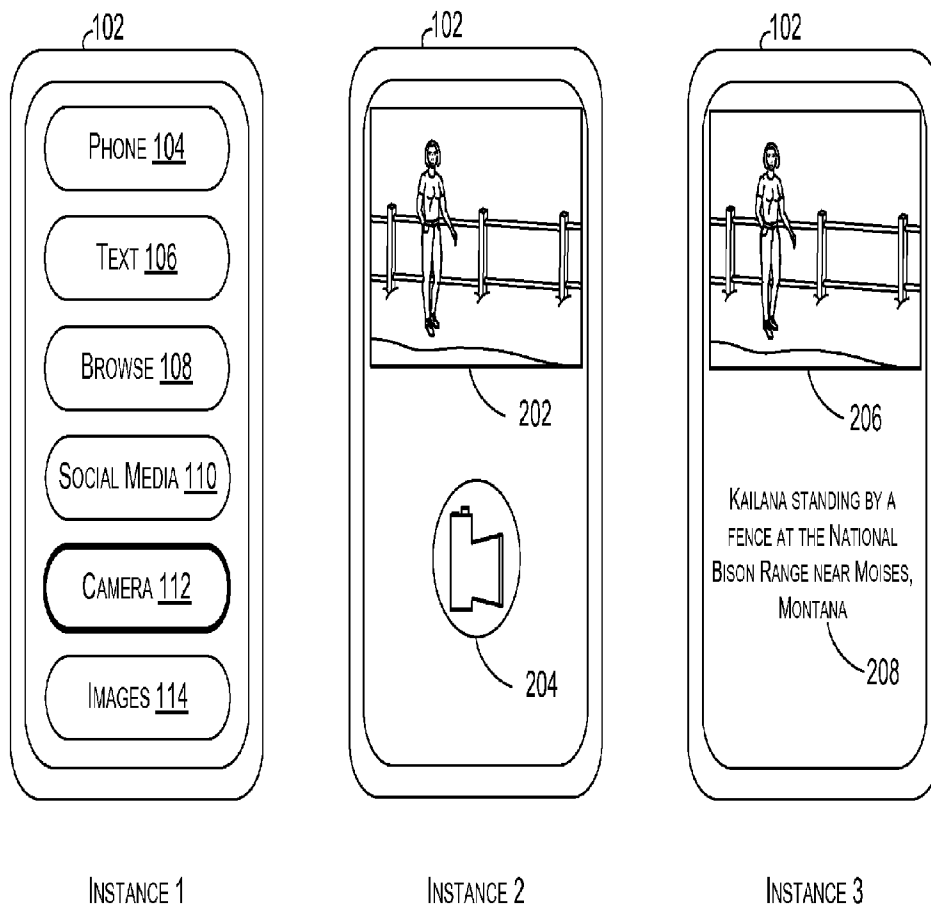


FIG. 2A

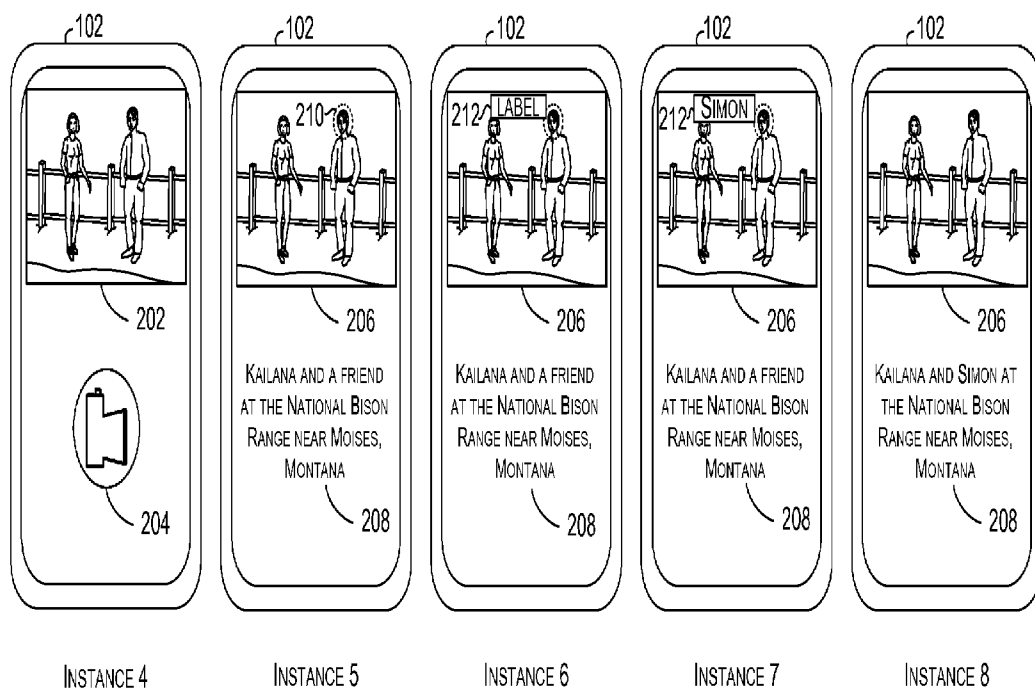


FIG. 2B

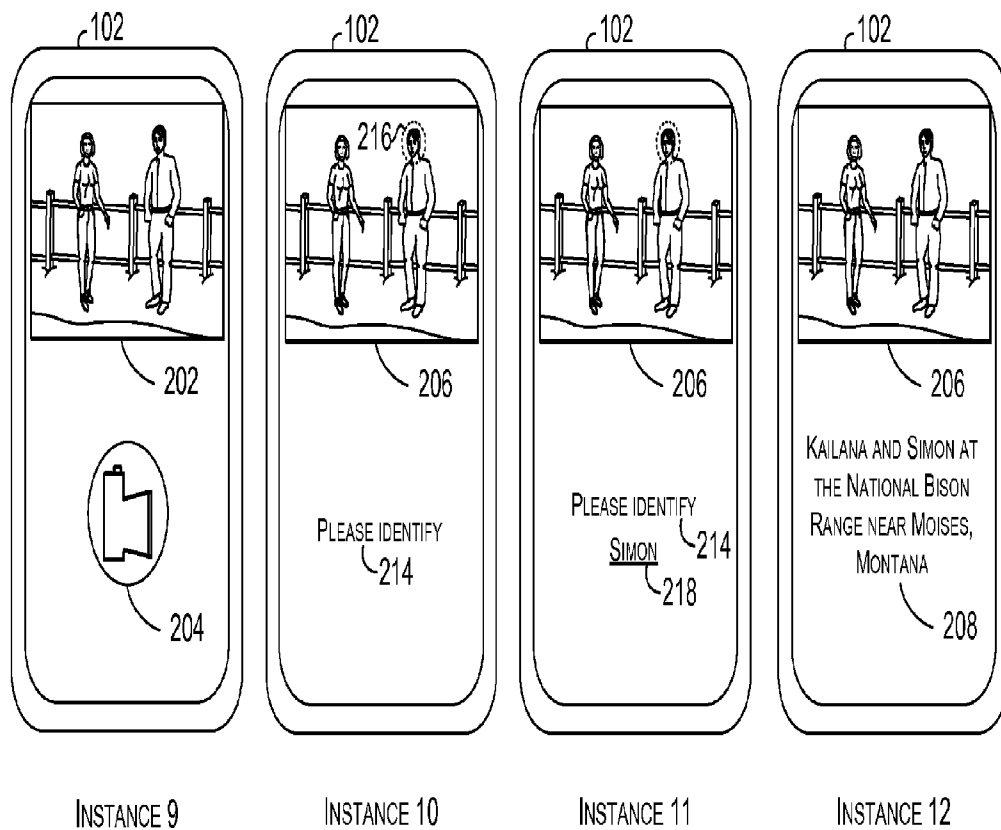


FIG. 2C

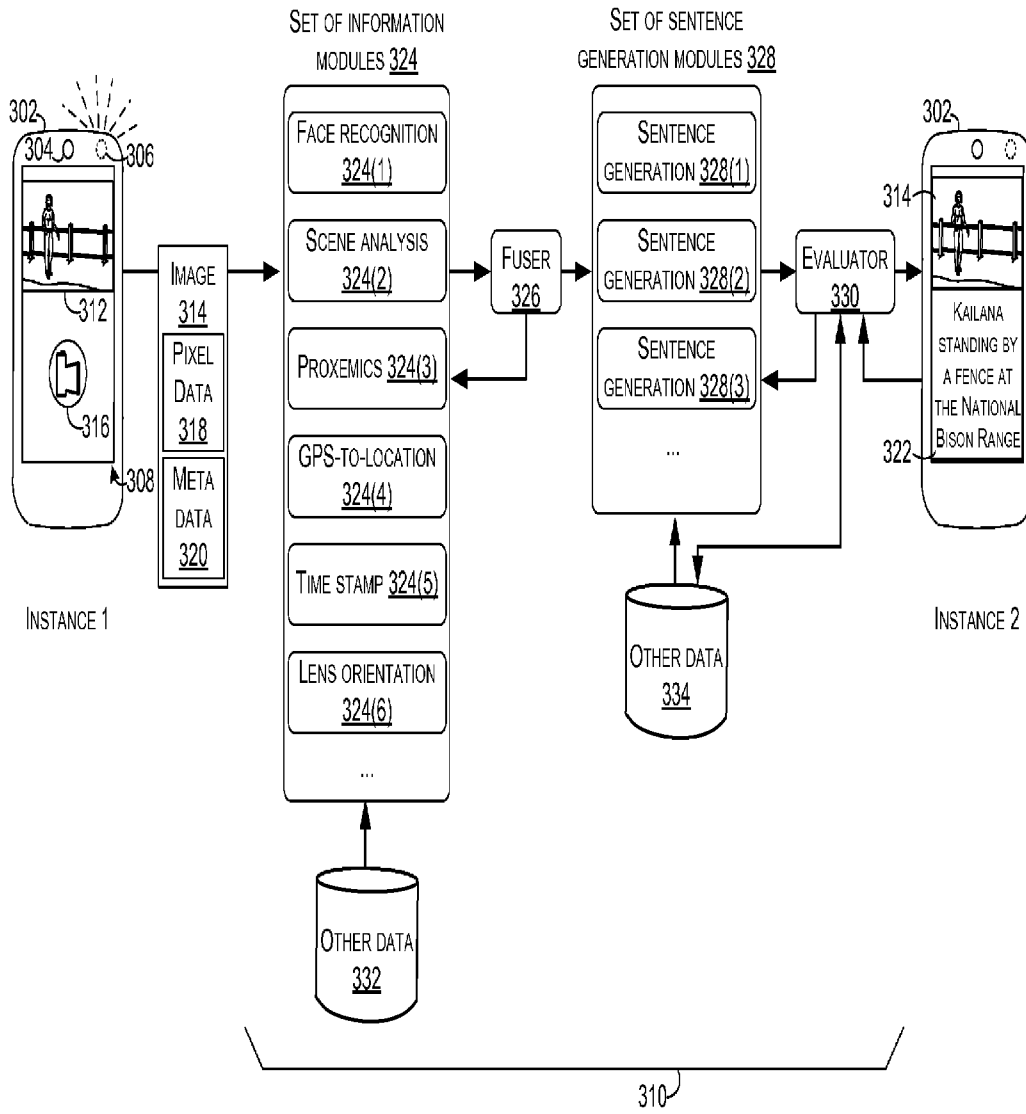


FIG. 3

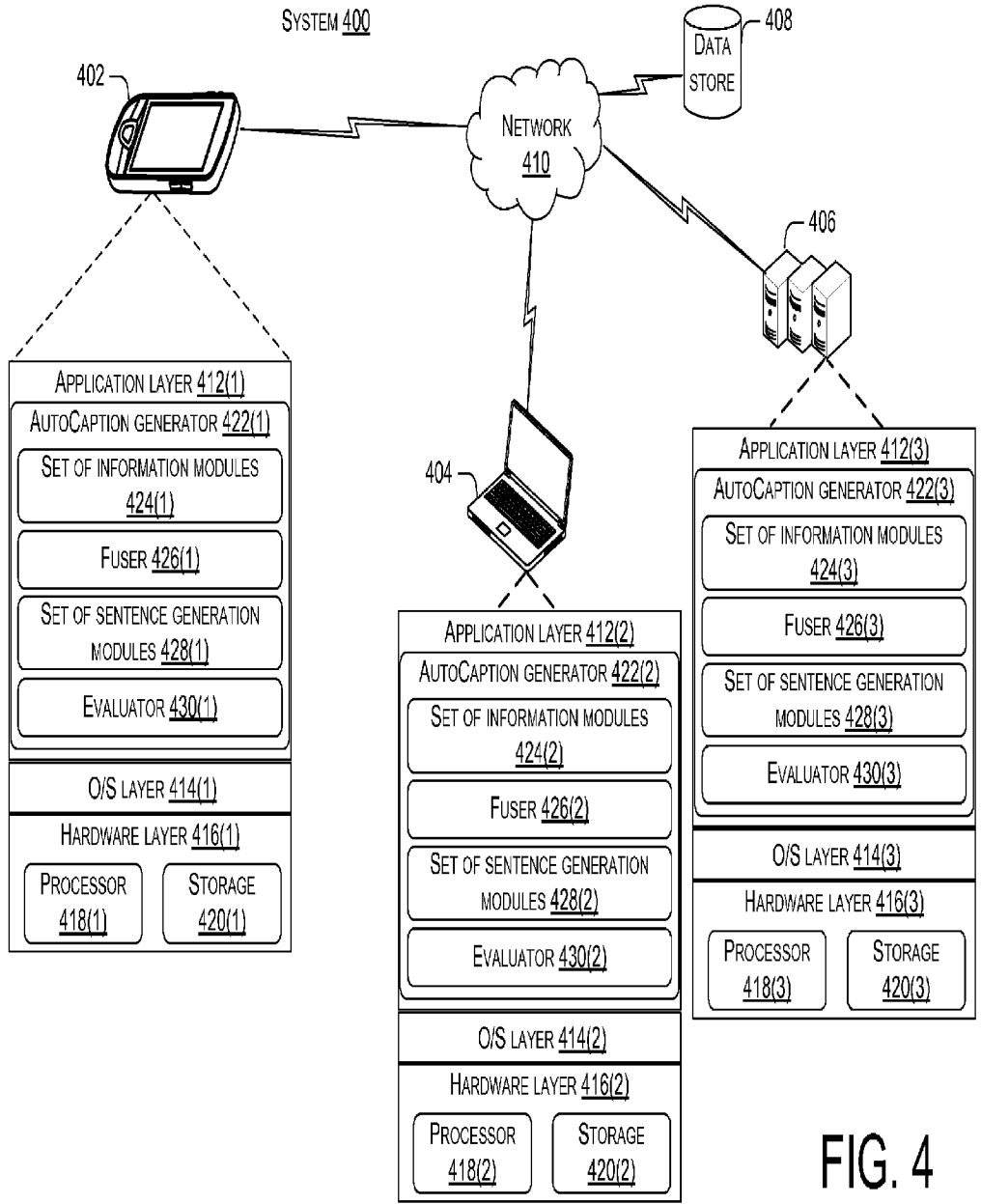


FIG. 4

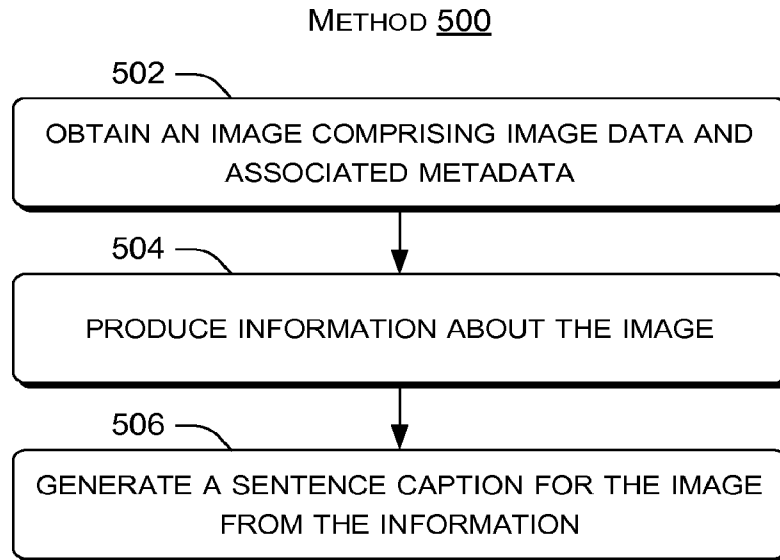


FIG. 5

AUTOCAPTIONING OF IMAGES

BACKGROUND

[0001] Digital photography has greatly increased the number of pictures that people (i.e., users) take. Users also enjoy the digital platform so that they can download (and/or upload) their pictures to a social media site to share with other people. Having a description (e.g., a caption) of the pictures enhances the viewing experience for the people with whom the pictures are shared. However, manually creating the captions is so time consuming that users rarely take the time to do it. Thus, the people who view the shared pictures often have a less than satisfying viewing experience.

SUMMARY

[0002] The described implementations relate to autocaptioning sentences to images. One example can include an image sensor, a processor, and a display. The image sensor can be configured to capture an image manifest as image data or pixel data. The processor can be configured to associate metadata with the image. The processor can also generate a sentence caption from the image data and the metadata. The display can be configured to present the image and the associated sentence caption. As used herein, a “sentence” for captioning an image can mean a sentence fragment, a complete sentence, and/or multiple sentences.

[0003] Another example can include a set of information modules and a set of sentence generation modules. The set of information modules can include individual information modules configured to operate on an image or metadata associated with the image to produce image information. The set of sentence generation modules can include individual sentence generation modules configured to operate on the image information to produce a sentence caption for the image.

[0004] The above listed examples are intended to provide a quick reference to aid the reader and are not intended to define the scope of the concepts described herein.

BRIEF DESCRIPTION OF THE DRAWINGS

[0005] The accompanying drawings illustrate implementations of the concepts conveyed in the present application. Features of the illustrated implementations can be more readily understood by reference to the following description taken in conjunction with the accompanying drawings. Like reference numbers in the various drawings are used wherever feasible to indicate like elements. Further, the left-most numeral of each reference number conveys the figure and associated discussion where the reference number is first introduced.

[0006] FIGS. 1A-1C and 2A-2C illustrate a computing device upon which sentence autocaptioning of images can be accomplished in accordance with some implementations of the present concepts.

[0007] FIG. 3 illustrates an example of a computing device upon which sentence autocaptioning of images can be accomplished and FIG. 3 also shows a process for accomplishing sentence autocaptioning of images in accordance with some implementations of the present concepts.

[0008] FIG. 4 illustrates a system which can accomplish sentence autocaptioning of images in accordance with some implementations of the present concepts.

[0009] FIG. 5 is a flowchart of a sentence autocaptioning process that can be accomplished in accordance with some implementations of the present concepts.

DETAILED DESCRIPTION

Overview

[0010] This patent relates to sentence autocaptioning of images. Pixel data and/or metadata of an image can be evaluated to obtain information about the image. This information can then be utilized to automatically generate a sentence autocaption for the image.

[0011] For purposes of explanation consider FIGS. 1A-1C which show a computing device 102 in several consecutive instances (e.g., FIG. 1A shows instances 1-4, FIG. 1B shows instances 5-7, and FIG. 10 shows instances 8-10). In this example, computing device 102 is manifest as a smart phone type mobile computing device that can provide multiple functionalities or features to the user. As illustrated in instance 1, the functionalities can include a phone functionality 104, a text message functionality 106, a web browsing functionality 108, a social media functionality 110, a camera functionality 112, and an images functionality 114. Assume for purposes of explanation that in instance 1, the user has selected the images functionality 114 as represented by bold highlight 116.

[0012] Instance 2 shows the computing device 102 responsive to the user selection related to instance 1. In this case, the computing device is now displaying options relating to the images functionality 114. For example, three folders or albums 120 (e.g., albums of images) are illustrated. The first album 120(1) is labeled “Christmas”, the second album 120(2) is labeled “Spring break”, and the third album 120(3) is labeled “Vacation”. The computing device also shows an “autocaption” function 122 that is discussed relative to instance 3. Assume that at instance 2 the user wants to view the images (e.g., pictures) in the vacation album and as such selects the vacation album as indicated by bold highlight 124.

[0013] Instance 3 shows three images 126(1), 126(2), and 126(3) of the vacation album 120(3). Assume that the user selects image 126(2) as indicated by the bold highlight of that element. Further, assume that the user wants a sentence caption to be automatically generated for the image and as such selects the “autocaption” function 122 as indicated by the bold highlight of that element.

[0014] Instance 4 is generated responsively to the user selections at instance 3. At 130, this example shows image 126(2) with an automatically generated sentence caption “Murdoch, Kailana, and Talise enjoying the beach in Poipu, Kauai.” Thus, sentence captions can be autogenerated for the user’s photos with little or no effort from the user.

[0015] FIG. 1B starts with instance 5 which is an alternative configuration to instance 4 of FIG. 1A. In this case, the user is also offered the choice to cause another sentence autocaption to be generated, such as for a case where the user does not like the generated sentence. In this example, the user choice is manifest as a ‘try again’ icon 132. Assume that the user selected the ‘try again’ icon 132 as indicated by the bold highlight of that element.

[0016] Instance 6 shows another autogenerated sentence caption “Murdoch, Kailana, and Talise walking along the surf in Poiou, Kauai.” at 130. The user can either select this sentence (e.g., associate the sentence as an autocaption) or cause another sentence to be generated via icon 132.

[0017] Instance 7 shows still another variation where multiple automatically generated sentences are presented for the user. The user can then select one of the sentences for auto-captioning of the image 126(2). In this example, two sentences are presented to the user at 130. The first sentence reads “Murdoch, Kailana, and Talise enjoying the beach in Poipu, Kauai.” The second sentence reads “Murdoch, Kailana, and Talise by a starfish.” The user can select either of the sentences. Alternatively, though not shown the user could be offered the ‘try again’ icon mentioned above to view other sentences.

[0018] FIG. 1C shows additional features of some implementations. Instance 8 is similar to instance 5 of FIG. 1B where the user selects image 126(2) and a sentence is automatically generated for the image at 130. In this case, the sentence reads “Murdoch, Kailana, and Talise by a starfish on the beach.” Instance 9 shows the user editing the sentence. In this case, the user is re-ordering the sentence by selecting and dragging (e.g., drag-and-drop) the word “starfish” to the beginning of the sentence as indicated by arrow 134. Instance 10 shows the result where the user editing is utilized as feedback and a new sentence is generated that reflects the user feedback. In this case, the updated sentence reads “A starfish on the beach with Murdoch, Kailana, and Talise.” Of course, other types of user edits are contemplated.

[0019] FIGS. 2A-2C collectively discuss autocaptioning features relative to a user taking a picture. These features can be interchangeable with the ‘album’ scenario described above relative to FIGS. 1A-1C and vice versa.

[0020] FIG. 2A begins with instance 1 which is similar to instance 1 of FIG. 1A except that the user has selected to use the camera functionality 112 to take a picture as indicated by the bold highlight of that element. Instance 2 shows an image preview 202. The user can take the picture by selecting the camera icon 204.

[0021] Instance 3 shows the resultant image (e.g., picture) 206 and associated auto-generated sentence at 208. In this example, the sentence reads “Kailana standing by a fence at the National Bison Range near Moises, Montana.”

[0022] Instance 4, of FIG. 2B, shows another image preview 202. Assume that the user presses the camera icon 204 to take the picture which is reflected in instance 5. Instance 5 shows a picture 206 and an autogenerated sentence at 208. In this example, the autogenerated sentence reads “Kailana and a friend at the National Bison Range near Moises, Montana.” Assume that the user wants to identify the other person (e.g., the friend) and as such taps the face on the screen as indicated at 210.

[0023] Instance 6 is generated responsive to the user action of instance 5 where the user tapped a portion of the image (in this case corresponding to a person’s face). A text box is shown at 212 into which the user can label the person. The user can enter text via a keyboard, a virtual keyboard, or other mechanism, such as by saying the person’s name. Any type of autosuggest and/or autocomplete technology can be utilized to make entry easier for the user.

[0024] Instance 7 shows the user entry label of “Simon” now associated with the image. This information can be utilized to update the autogenerated sentence. Instance 8 shows the results where the updated autogenerated sentence indicated at 208 now reads “Kailana and Simon at the National Bison Range near Moises, Montana.” This feature is not limited to the user labeling (e.g., tagging) faces. For instance, the user could label or tag a building as “my favorite restaur-

ant” our “our beach house.” The user could alternatively label the natural environment, such as “our secret beach.”

[0025] FIG. 2C shows another variation. In this case, instance 9 shows the image preview and instance 10 shows the image 206. However, in instance 10 a request 214 is made to the user to identify the circled person 216. In instance 11, the user enters “Simon” at 218. This user-provided information can then be utilized to automatically generate sentence 208 which reads “Kailana and Simon at the National Bison Range near Moises, Montana.” Thus, techniques are described to allow user labeling of pictures on their own initiative or when requested. This information can be utilized to enhance the quality of the autogenerated sentences, for that image and/or other images.

[0026] In the above discussion single sentences are automatically captioned to the image. However, other implementations can generate sentence fragments, single sentences, and/or multiple sentences for an image. In this case, multiple sentences means sentences which augment one another, such as sentences of a paragraph, rather than alternative sentences. For instance, in the example discussed relative to instance 12, the sentences might read “Kailana and Simon at the National Bison Range near Moises, Montana. They are standing near the fence that borders the Western edge of the Range.”

[0027] An example process for automatically generating sentence captions (e.g., sentence autocaption) is explained below relative to FIG. 3. Note that the illustrated computing device configuration is provided for purposes of explanation and other configurations are contemplated that allow sentence autocaptions to be generated for images.

[0028] FIG. 3 shows another computing device 302. In this case, the computing device includes two cameras; a first camera 304 points toward the user (e.g., out of the drawing page) and a second camera 306 points away from the user (e.g., into the drawing page). The second camera is shown in dotted lines (e.g., ghost) since it would not actually be visible to the user from the front of the computing device 302. Computing device 302 also includes a touch sensitive screen or display 308.

[0029] For purposes of explanation, computing device 302 is shown at instance 1 and instance 2. Instance 1 and instance 2 are separated by process 310 which is described below.

[0030] Beginning with instance 1, assume that the user has selected the camera functionality associated with camera 306. As a result, an image preview 312 is presented on display 308 and that the user takes a picture or image 314 by selecting icon 316. The image 314 can include image or pixel data 318 and associated metadata 320. The metadata can convey GPS coordinates that the image was captured at, date and time of image capture, type of camera, and/or type of computing device, among other information. As seen in instance 2, process 310 can produce a sentence autocaption 322 that is associated with the image 314.

[0031] In this example, process 310 can be accomplished by a set of information modules 324, an information fuser (or “fuser”) 326, a set of sentence generation modules 328, and a sentence evaluator 330. In this case, the set of information modules can include a face recognition module 324(1), a scene analysis module 324(2), a proxemics module 324(3), a GPS-to-location module 324(4), a time stamp module 324(5), and a lens orientation module 324(6), among others. Other information modules can alternatively or additionally be included. Further, multiple instances of a given module type

can be employed. For example, two or more face recognition modules could be employed instead of just one.

[0032] In this case, the set of information modules **324** is shown in parallel receiving a single input (e.g., the image **314**) and generating a single output (to fuser **226**). However, such need not be the case. For example, individual information modules can receive the input and generate their own output. Further, the modules can be arranged in a serial fashion such that the output of one module can serve as input to another module. One configuration can be thought of as a pipeline configuration. In such a configuration, several information modules can be arranged in a manner such that the image and output from one module serve as input to another module to achieve results that the another module cannot obtain operating on the image alone. Further still, even in a parallel relationship, one module can utilize the output of another module. One such example is described below relative to the face recognition module **324(1)**, the scene analysis module **324(2)** and the proxemics module **324(3)**.

[0033] Briefly, individual information modules **324** can process the image's pixel data **318** and/or metadata **320** and/or other data **332** to produce information that is potentially useful to the set of sentence generation modules **328**. The other data **332** can include many different types of data. For instance, the other data **332** can include faces that were previously labeled by the user relative to this or other images. (Such an example is described above relative to the user labeling or tagging a face in the image). In another case, the other data **332** could include the time stamps of other images that are already associated with an autocaption sentence. For instance, if the present image was taken in close temporal relation to this one (e.g., less than a second apart) the same elements may be in both images. As such the already labeled elements of the earlier captioned image can be useful for analyzing the present image. This aspect can be especially valuable where the image is from a set of images captured in a camera "burst mode" or where the image is a video frame from a video. In summary, images can be related in various ways. For instance, the images can be related temporally, by location, by objects or people in the images, etc. The relation can be part of other data **332** and can be utilized to derive useful information about individual images.

[0034] In some cases, the facial recognition module **324(1)** can operate cooperatively with a face detection module (not shown for sake of brevity). The face detection module can identify regions in the pixel data **318** that are likely to be human faces. The facial recognition module **324(1)** can then further process these regions to attempt to identify the face. For example, the facial recognition module can analyze various parameters of the face and compare values of those parameters to parameter values of labeled faces from other data **332**. In this case, the facial recognition module **324(1)** can receive an indication that one face occurs in the image and identify the face as belonging to Kailana. This information can be the output from the face recognition module to the fuser **326**.

[0035] The scene analysis module **324(2)** can identify various features or elements of the image **314**. In this case, the scene analysis module can identify the vertical and horizontal features of the fence and identify these features as a fence. The scene analysis module can also receive output from the face recognition module regarding the location and identity of a face. The scene analysis module can use this information to

identify the location of Kailana's body on the image **314**. This information can be output to the fuser **326**.

[0036] In another implementation, the scene analysis module **324(2)** can function to determine the scene or scenario of the image. For instance, the scene analysis module can determine whether the image was taken: (a) at the beach, (b) in a park, (c) in a garden, (d) while skiing, etc. In such a case, the emphasis of the scene analysis module can be to identify the scene or scenario, rather than identifying specific objects within the image.

[0037] The proxemics module **324(3)** can serve to identify relative relationships between elements of the image **314**. Often the proxemics module functions to determine the relationship between two humans, such as whether the humans are embracing, holding hands, shaking hands, etc. In this case, the proxemics module can receive input from the scene analysis module **324(2)** regarding the location of Kailana and the location of the fence and determine their orientation. For instance, in this example Kailana is near the fence, but is not engaging the fence, such as jumping over the fence or sitting on the fence. The proxemics module **324(3)** can output this information to the fuser **326**.

[0038] The GPS-to-location module **324(4)** can operate on the image's metadata **320** to convert GPS coordinates of the metadata into a textual location. In this case, the GPS-to-location module maps the GPS coordinates to the National Bison Range near Moises, Montana., and outputs this information to the fuser **326**.

[0039] The time stamp module **324(5)** can also operate on the metadata **320**. The time stamp module can analyze the timestamp to determine useful information such as the month or season, the day, the day of the week, and/or the time of day. For example, it can be useful for sentence generation to know whether the picture was taken in the Spring, Summer, Fall, or Winter. Also, the fact that the picture was taken on a particular day, such as 4th of July or Christmas, can be useful for sentence generation. Further, the time of the day can be useful (e.g., whether the photo was taken in the morning or evening).

[0040] In this case, the time stamp module **324(5)** can also compare the time stamp from the metadata to other time stamps in the other data **332** to determine if any other images were taken in close temporal relation to this image. If so, the time stamp module can provide information from the other images to the fuser **326** and/or to other information modules **324**. For instance, if the computing device **302** took another picture a half second before this one that contained one face that was already identified as belonging to Kailana, this information can be utilized by the face recognition module to increase its confidence score that it has correctly identified the face in this image (e.g., image **314**). This type of scenario is especially useful where the image is from a "burst mode" or where the image is part of a video stream.

[0041] The lens orientation module **324(6)** can evaluate the metadata **320** to determine which camera took the image **314**. For instance, in cases such as this one where the backwards facing camera **306** is used, then the person taking the picture (e.g., the user) is not likely to be in the image **314**. In contrast, if the forward facing camera **304** is utilized, the user is likely taking a picture of themselves in front of an interesting scene. This information can be sent to the fuser **326** for utilization by the set of sentence generation modules **328**. For example, the sentence "Kailana taking a picture of herself at the entrance to Yellowstone" may be more interesting than "Kailana at the

entrance to Yellowstone.” This aspect will be discussed in more detail below relative to the set of sentence generation modules 328.

[0042] Of course, the listed information modules of the set of information modules 324 are provided for purposes of explanation. Other examples of information modules that can be employed can include weather recognition modules, landmark recognition modules, celebrity recognition modules, animal recognition modules, facial expression modules (e.g. whether a person is smiling or frowning), activity recognition modules (e.g., “having dinner”, “playing tennis”, “skiing”), clothing recognition modules (e.g., wearing a red shirt) and/or car model recognition modules, among others. Further, the present concepts allow the set of information modules to be readily configurable to accept new and/or different information modules and/or for existing information modules to be removed, as desired.

[0043] The fuser 326 can receive the potentially useful information produced by the individual information modules of the set of information modules 324 about the image 314. The fuser 326 can process, evaluate, filter, and/or organize the information for receipt by the set of sentence generation modules 328. For instance, in an implementation employing two face recognition modules 324(1), the first module could identify a face in the image 314 as one person, whereas the second module can identify the same face as another person or determine that there is not a face in the image. In such a scenario, the fuser can determine what information to forward to the set of sentence generation modules 328. For example, the fuser could decide that one module is more reliable than the other and only pass on its (the more reliable module’s) determination. In another case, the fuser could pass on both findings along with their relative probabilities to the set of sentence generation modules 328.

[0044] In still another example, the fuser 326 can receive potentially incompatible information from two different information modules. For instance, the GPS-to-location module 324(4) might indicate that the picture was taken in Denver, Colo. and the scene analysis module 324(2) might indicate that the picture was taken at the beach. In such a case, the fuser may weight the reliability of the information from the GPS-to-location module higher than the information from the scene analysis module since beaches tend to occur with less frequency in Colorado than in states that border an ocean. Further, the fuser may request the scene analysis module to repeat its analysis to determine if the image was mis-identified. For instance, maybe snow was mis-identified as sand, for example. If the fuser is confident in the information from both modules, the fuser can also provide this information in a manner that allows for the generation of a more accurate sentence. For instance, many people tend to associate the beach with the ocean. So in this case, a resultant sentence might be improved by saying “. . . on a South Platte River Beach near Denver, Colorado.”

[0045] In summary, the fuser 326 can function to provide a filtered set of germane and accurate information to the set of sentence generation modules 328 and thereby enhance the quality of the sentences produced by the set of sentence generation modules 328. Alternatively or additionally, this filtering can reduce the processing resources and/or time consumed by the set of sentence generation modules 328 to generate a sentence from the information. Note also, that the fuser may not operate only in a one-way relationship between the set of information modules 324 and the set of sentence

generation modules 328. Instead, the fuser 326 can feed information back to the set of information modules 324 in an iterative fashion to allow the set of information modules 324 to increase the quality of (e.g., refine) information that they output. The fuser can also use this feedback information to improve its own performance. That said, some implementations can operate without a fuser 326 and instead image related information output by the set of information modules 324 can be input directly into the set of sentence generation modules 328.

[0046] The set of sentence generation modules 328 can receive the information from the set of information modules 324 directly or a sub-set from the fuser 326. In this case, the set of sentence generation modules is represented as sentence generation modules 328(1)-328(3) but any number of one or more sentence generation modules can be utilized in various implementations. The set of sentence generation modules 328 can utilize this information from the fuser 326, as well as information from other data 334 to generate sentence fragments (with or without punctuation) and/or sentences for the image 314. For example, the other data 334 can include a training corpus, previously generated auto-caption sentences, user feedback, user style preferences, etc.

[0047] The other data 334 can allow the set of sentence generation modules 328 to generate sentence captions that are objectively correct and configured to please the subjective preference of the user. In a relatively simple example, the other data 334 can help the sentence generation modules 328 to avoid generating repetitive sentences. For instance, assume that the user quickly takes two pictures of Murdoch at the beach. Relative to the first picture, the sentence generation modules may autogenerate the sentence “Murdoch at the beach.” This sentence can be added to the other data 334 so that the sentence generation modules are less likely to generate the same sentence for the second picture. Instead the generated sentence for the second image might read “Murdoch at the edge of the ocean” or “Another shot of Murdoch at the beach.” While all three of these sentences are objectively accurate, a second occurrence of the first sentence may be met with less enthusiasm from the user than the latter sentences.

[0048] In another example, the other data 334 may indicate that the user likes more specificity in the autogenerated sentence captions 322. For instance, in a previous session, assume that an autogenerated sentence caption read “John in the mountains” and that the user had edited the caption to read “John in the Rocky Mountains.” This user feedback information contained in the other data 334 can be applied to the above example so that rather than autogenerating a sentence caption like “Murdoch at the beach” the sentence generation modules can autogenerate a sentence caption like “Murdoch at Poipu Beach.” The individual sentence generation modules can generate sentences from the information obtained from the fuser 326 and the other data 334. These sentences can be processed by the evaluator 330.

[0049] In one implementation, individual sentence generation modules 328 can employ sentence templates to automatically generate the sentences for the images. One example of a template can be “[Person1] and [Person2] at the [Scene] in [City Location].” For example, assume that the face recognition module 324(1) finds Person1=Kailana and Person2=Talise and scene analysis module 324(2) returns Beach and GPS-to-location module 324(4) returns “Poipu, Kauai.” Then sentence generation module 328 can generate “Kailana and Talise at the beach in Poipu, Kauai.” Some

implementations can maintain a large number of templates, each covering a different subset of possible results of the sentence generation modules. The evaluator 330 and/or the sentence generation modules 328 can determine which of the templates are applicable and choose one (such as the one where the most blanks are filled in).

[0050] The evaluator 330 can evaluate the sentences obtained from the set of sentence generation modules 328 to select a sentence to autocaption with the image 314. The evaluator can utilize information from other data 334 in the selection process and/or direct user feedback. For instance, if the evaluator receives one sentence from each sentence generation module 328(1), 328(2), and 328(3) the evaluator can compare the sentences to one another relative to one or more parameters to rank the sentences. For instance, the parameters can relate to sentence style (e.g., writing style), sentence detail and/or repetitiveness, user preferences, and/or user feedback, among others.

[0051] In some cases, the choice made by the evaluator 330 can be weighted by the parameters such as user preferences and/or previously generated sentences. For instance, the user may indicate that he/she prefers one template style over another template style. In such a case, this user preference can be used as a factor in selecting or choosing from the available templates. Thus, the type of template favored by the user can be a positive parameter. Similarly, if one template has been chosen recently to generate a sentence, that template can be a negative parameter to reduce the chance that identical (or very similar sentences) are selected for related pictures (e.g., those pictures taken in temporal sequence or occurring in the same album, for example).

[0052] Template style (e.g., writing style) can be proactively user defined or inferred from user feedback. For example, when autocaptioning is first set-up or activated the user could be presented with a series of writing samples and the user can select his/her preferences as a guide for future autocaption sentence structure. In another example, the evaluator 330 may iteratively utilize user feedback of autocaptioned sentences to develop a sentence style for the user.

[0053] As mentioned above, sentence detail can relate to how much detail is included in a sentence and/or what the detail relates to. For example, the user may want geographic details included in the autocaption sentences (e.g., “. . . at the beach in Poipu, Kauai” rather than “. . . at the beach”). The detail may also relate to time, naming, etc. For instance, the user may prefer “. . . walk along the beach at sunrise” rather than “. . . walk along the beach.” Similarly, the user may want full names rather than less formal names. For example, the user may prefer “John Smith walks on the beach” to “John walks on the beach” (or vice versa). The evaluator 330 can consider (and in some cases weight) these various parameters to rank the sentences obtained from the set of sentence generation modules 328. The evaluator can select the highest ranking sentence. Alternatively, the evaluator can compare the sentences to a threshold and present any sentences that satisfy the threshold to the user so that the user can make the final selection.

[0054] Stated another way, the evaluator 330 can select a single sentence or multiple sentences to autocaption (as indicated at 322) for the image 314. For instance, the evaluator 330 may initially select two sentences from the set of sentence generation modules 328 to autocaption the image. If the user selects one of the sentences, this direct user feedback can be utilized for selecting a sentence for autocaptioning subse-

quent images. Despite the potential advantages offered by the evaluator 330, some implementations, such as those employing a single sentence generation module in the set of sentence generation modules 328, can function without the evaluator.

[0055] Note that the process 310 can serve to accomplish the sentence autocaptioning functionality. The process 310 can occur completely on the computing device 302. In other implementations, some or all of the process can be performed external to the computing device, such as by server or cloud-based resources.

[0056] Note also that a potential advantage of the described configuration is the flexible or configurable nature of the set of information modules 324 and/or the set of sentence generation modules 328. In this configuration, new information modules can be added to the set of information modules 324 as they become available and/or when existing modules are improved (e.g., upgraded). Similarly, new sentence generation modules can be added to the set of sentence generation modules 328 as they become available and/or when existing modules are improved (e.g., upgraded). The additions and/or deletions to the modules can be accomplished in a seamless fashion that need not interfere with providing sentence autocaptioning of images to the user.

[0057] Note further, that while process 310 is explained relative to a single image, the process lends itself to use with sets of images. In one case, a set of images can be an album. The feedback mechanisms provided in process 310 can allow the images of an album to be labeled with autocaption sentences that are interesting to read, non-repetitive and contextually related. Stated another way, the autocaption sentences can be contextually related in such a way that taken as a whole, the autocaption sentences provide a narrative for the album. For instance, if the album relates to a vacation, the first image may be autocaptioned with a sentence like “Our vacation started early in the morning at the San Jose, California Airport.” The next autocaption sentence might state that “We landed in Kauai in the afternoon.” The next autocaption sentence might read “On the beach at last.” As such, the autocaption sentences can collectively convey the essence of the images of the album. Further the autocaption sentences can be utilized in various ways. For instance, the autocaption sentences can be read aloud by a text-to-audio component to aide a visually impaired user in appreciating the contents of the album.

[0058] Further still, concepts can be applied to a set of images taken in close temporal proximity to provide meaning to the set. For instance, a burst mode of images can be analyzed so that the autocaption sentences that are generated for individual images of the set provide a narrative to the set. For instance a set of images showing a baseball batter could be labeled as “John awaits the pitch,” “John swings at the pitch,” and “John strikes out.” The same technique can be applied to video. For instance, one implementation could autogenerate a sentence caption for frames occurring every second and/or every time there is a scene change. These autogenerated sentence captions can describe the video as a whole. The autogenerated sentence captions can be utilized for various purposes. For instance the autogenerated sentence captions can be searched to find specific events in the video. For instance, if the user wants to view the portion of the video that was taken at the National Bison Range, the user can select to search the autocaption sentences for this phrase and the corresponding portions of the video can be displayed for the user.

[0059] In summary, the above described implementation can take an image and then run a number of different algorithms that analyze various aspects of the image (image content, EXIF (exchangable image file format) data, user metadata, etc.) and automatically generate a sentence caption for the image. From one perspective, the present concepts can leverage the relationships between images to generate accurate and pleasing autocaption sentences. For example, the images may share the relation of being in the same album or folder, taken together in a burst, taken on the same day, taken at the same location, contain the same faces, etc. to generate meaningful autocaptions from sets of modules. The modules included in the sets can be easily changed and/or updated to refine the results.

System Example

[0060] FIG. 4 shows a system 400 configured to accomplish sentence autocaptioning of images. In this case, system 400 includes three computing devices 402, 404, and 406, as well as a data store 408. These components can be coupled via a network 410. For purposes of explanation, assume that computing device 402 (similar to computing device 102 of FIGS. 1A-1C, computing device 202 of FIGS. 2A-2C, and computing device 302 of FIG. 3) is a user's mobile device, such as a smart phone type device and that computing device 404 is the user's notebook computer. Also assume that computing device 406 is a remote computing device or devices that can provide computing resources. For instance, computing device 406 can be manifest as a server computer or as cloud-based computing resources.

[0061] Computing devices 402, 404, and 406 can include several elements which are defined below. For example, these computing devices can include an application layer 412 that operates upon an operating system layer 414 that operates upon a hardware layer 416. (In this case, the suffix "(1)" is used to indicate an instance of these elements on computing device 402, while the suffix "(2)" is used to indicate an instance on computing device 404 and the suffix "(3)" is used to indicate an instance on computing device 406. The use of these elements without a suffix is intended to be generic).

[0062] The hardware layer 416 can include a processor 418 and storage 420, as well as additional hardware components, such as input/output devices, buses, graphics cards, cameras, lenses, image sensors, etc., which are not illustrated or discussed here for sake of brevity.

[0063] The application layer 412 can include an autocaption generator 422. The autocaption generator 422 can include a set of information modules 424, a fuser 426, a set of sentence generation modules 428, and an evaluator 430. (These elements are similar to, or the same as, the like named elements introduced relative to FIG. 3).

[0064] The set of information modules 424 can include individual information modules configured to operate on an image and/or metadata associated with the image to produce image information. Information modules can be readily added to, or removed from, the set of information modules.

[0065] The set of sentence generation modules 428 can include individual sentence generation modules configured to operate on the image information to produce a sentence caption for the image. Sentence generation modules can be readily added to, or removed from, the set of information modules.

[0066] System 400 enables various implementations that can generate sentence captions for images based on the analy-

sis of image content and related metadata. The autocaption generator 422 can allow users to capture images and display automatically generated captions for those images. These images and captions can then be shared via a number of social networking media, through the autocaption generator 422.

[0067] In one implementation, all three of computing devices 402, 404, and 406 can be configured to accomplish the sentence autocaptioning of image concepts described above and below. For example, computing device 402 can have a robust autocaption generator 422(1), such that computing device 402, operating in isolation, can accomplish sentence autocaptioning of images. Similarly, computing device 404 can have a robust autocaption generator 422(2), such that computing device 404, operating in isolation, can accomplish sentence autocaptioning of images.

[0068] In other implementations, the computing device 406 can have a relatively robust autocaption generator 422(3), while computing device 402 may include an autocaption generator 422(1) that offers more limited functionality for accomplishing the described concepts in cooperation with the computing device 406. In one such case, the autocaption generator 422(1) on computing device 402 can include a communication component (not shown). The communication component can be configured to obtain an image captured on the computing device 402, send the image to computing device 406 which can perform a majority of the data storage and processing. For example, the autocaption generator 422(3) of computing device 406 can generate the sentence autocaption for the image and return the sentence autocaption to computing device 402. The autocaption generator 422(1) of computing device 402 can associate the sentence autocaption with the image for presentation to the user and/or storage.

[0069] Other configurations can alternatively or additionally be employed. For example, computing device 406 alternatively or additionally to sending the sentence autocaption to computing device 402, could send the image and the autocaption to the user's notebook computing device 404 for storage and/or subsequent use by the user. In still another configuration, computing device 406 may or may not return the autocaption to computing device 402, and can send the autocaption and the image to datastore 408. Datastore 408 can then act as a remote repository (e.g., cloud-based) of the user's autocaptioned images. Alternatively or additionally, datastore 408 could serve to store the other data (332 and/or 334 of FIG. 3).

[0070] Stated another way, the autocaption generators 422 of system 400 can allow users to capture images and upload to a backend system (e.g., computing device 406). The backend system can automatically generate corresponding sentence captions that can then be displayed within an application (e.g., the autocaption generator 422(1)) on computing device 402. The user can then upload the images and captions to a social networking site (e.g., social media functionality 110 of FIG. 1) or share the autocaptioned images via different media. Affordances, such as viewing upload history and the ability to easily edit the auto generated captions, can also be enabled within the autocaption generators 422.

[0071] The term "computer" or "computing device" as used herein can mean any type of device that has some amount of processing capability and/or storage capability. Processing capability can be provided by one or more processors (such as processor 418) that can execute data in the form of computer-readable instructions to provide a functionality. Data, such as computer-readable instructions, can be

stored on storage, such as storage 420 that can be internal or external to the computer. The storage can include any one or more of volatile or non-volatile memory, hard drives, flash storage devices, and/or optical storage devices (e.g., CDs, DVDs etc.), among others. As used herein, the term “computer-readable media” can include signals. In contrast, the term “computer-readable storage media” excludes signals. Computer-readable storage media includes “computer-readable storage devices.” Examples of computer-readable storage devices include volatile storage media, such as RAM, and non-volatile storage media, such as hard drives, optical discs, and flash memory, among others.

[0072] In the illustrated implementation computing devices 402, 404, and 406 are configured with a general purpose processor 418 and storage 420. In some configurations, a computer can include a system on a chip (SOC) type design. In such a case, functionality provided by the computer can be integrated on a single SOC or multiple coupled SOCs. In one such example, the computer can include shared resources and dedicated resources. An interface(s) can facilitate communication between the shared resources and the dedicated resources. As the name implies, dedicated resources can be thought of as including individual portions that are dedicated to achieving specific functionalities. For instance, in this example, the dedicated resources can include the autocaption generator 422.

[0073] Shared resources can be storage, processing units, etc. that can be used by multiple functionalities. In this example, the shared resources can include the processor. In one case, as mentioned above autocaption generator 422 can be implemented as dedicated resources. In other configurations, this component can be implemented on the shared resources and/or the processor can be implemented on the dedicated resources. In some configurations, the autocaption generator 422 can be installed during manufacture of the computer or by an intermediary that prepares the computer for sale to the end user. In other instances, the end user may install the autocaption generator 422, such as in the form of a downloadable application.

[0074] Examples of computing devices can include traditional computing devices, such as personal computers, cell phones, smart phones, personal digital assistants, pad type computers, cameras, or any of a myriad of ever-evolving or yet to be developed types of computing devices. Further, aspects of system 400 can be manifest on a single computing device or distributed over multiple computing devices.

First Method Example

[0075] FIG. 5 shows a flowchart of a sentence autocaptioning of images method or technique 500 that is consistent with at least some implementations of the present concepts.

[0076] At block 502 the method can obtain an image comprising image data and associated metadata.

[0077] At block 504 the method can produce information about the image.

[0078] At block 506 the method can generate a sentence caption for the image from the information.

[0079] Method 500 can be performed by the computing devices described above relative to FIGS. 1A-1C, 2A-2C, 3, and/or 4, and/or by other devices and/or systems. The order in which the method 500 is described is not intended to be construed as a limitation, and any number of the described blocks can be combined in any order to implement the method, or an alternate method. Furthermore, the method can

be implemented in any suitable hardware, software, firmware, or combination thereof, such that a computing device can implement the method. In one case, the method is stored on computer-readable storage media as a set of instructions such that execution by a computing device causes the computing device to perform the method.

CONCLUSION

[0080] Although techniques, methods, devices, systems, etc., pertaining to sentence autocaptioning of images are described in language specific to structural features and/or methodological acts, it is to be understood that the subject matter defined in the appended claims is not necessarily limited to the specific features or acts described. Rather, the specific features and acts are disclosed as exemplary forms of implementing the claimed methods, devices, systems, etc.

1-20. (canceled)

21. A device, comprising:

a display;

a processor; and,

storage storing computer-readable instructions which, when executed by the processor, cause the processor to: obtain an image comprising image data and associated metadata,

receive a label from a user, the label corresponding to an individual non-human element that is visible in the image,

access previous sentence captions associated with the user,

automatically generate sentence captions for the image based at least in part on the image data, the associated metadata, the label, and the previous sentence captions, and,

present the sentence captions on the display.

22. The device of claim 21, wherein the sentence captions comprise full sentence captions or sentence fragment captions.

23. The device of claim 21, wherein the computer-readable instructions further cause the processor to develop a sentence template style based at least in part on previous sentence captions associated with the user.

24. The device of claim 23, wherein the sentence template style further comprises preferences of the user.

25. The device of claim 21, wherein the computer-readable instructions further cause the processor to access a sentence template style of the previous sentence captions.

26. The device of claim 25, wherein the computer-readable instructions further cause the processor to automatically generate the sentence captions for the image using the sentence template style as a positive parameter.

27. The device of claim 25, wherein the computer-readable instructions further cause the processor to automatically generate the sentence captions for the image using the sentence template style as a negative parameter.

28. The device of claim 25, wherein the sentence template style is selected by the user.

29. The device of claim 21, wherein the computer-readable instructions further cause the processor to automatically generate sentence captions for the image based at least in part on the image being included in a set of images.

30. The device of claim 29, wherein the previous sentence captions are associated with other images in the set and the sentence captions for the image are contextually related to the previous sentence captions.

31. A method, comprising:
capturing an image with an image sensor, the image comprising image data;
associating metadata with the image;
based at least in part on the image data and the metadata, identifying a scenario involving a human and an object of the image;
based at least in part on the scenario, identifying a relative relationship between the human and the object in the image;
producing multiple sentence captions for the image that reflect the relative relationship between the human and the object; and
presenting a display of the multiple sentence captions.

32. The method of claim **31**, wherein the method is performed by a device that includes the image sensor and presents the display to a user.

33. The method of claim **31**, wherein identifying the scenario further comprises identifying a scene that includes the human and the object, and wherein identifying the relative relationship further comprises identifying an orientation of the human relative to the object.

34. A device, comprising:
an image sensor configured to capture images comprising pixel data;
a processor; and,
storage storing computer-readable instructions which, when executed by the processor, cause the processor to: generate image information from the pixel data of an individual image captured by the image sensor and metadata associated with the image, the image containing a human and a non-human element,

obtain a relative relationship between the human and the non-human element in the image,
access a template style of a user, and
automatically generate sentence captions for the image based at least in part on the pixel data, the metadata, the relative relationship, and the template style of the user.

35. The device of claim **34**, wherein the computer-readable instructions further cause the processor to determine the relative relationship between the human and the non-human element in the image.

36. The device of claim **34**, wherein the computer-readable instructions further cause the processor to obtain the relative relationship from another device.

37. The device of claim **34**, wherein the computer-readable instructions further cause the processor to send the pixel data and the metadata to another device for processing to determine the relative relationship.

38. The device of claim **37**, wherein the another device comprises cloud-based computing resources.

39. The device of claim **34**, wherein the image is associated with a set of related images, and wherein the computer-readable instructions further cause the processor to automatically generate the sentence captions by considering other sentence captions for other images in the set of related images.

40. The device of claim **39**, wherein the set of related images comprises an album and the sentence captions and the other sentence captions comprise a contextually related narrative of the album.

* * * * *