**SOFTWARE ARTICLE**                                                      **Open Access**

CrossMark

# A platform for creating Smartphone apps to enhance Chinese learning using augmented reality

Yiqun Li[*], Aiyuan Guo, Ching Ling Chin and Joo-Hwee Lim

## Abstract

We developed a software system consisting of a container app for smartphones and a web portal. The container app, named "MimasAR" (iOS container app: https://itunes.apple.com/sg/app/mimasar/id988692600?mt=8; Android container app: https://play.google.com/store/apps/details?id=sg.edu.astar.i2r.mimasar&hl=en) can be used to recognize Chinese characters on printed media quickly and reliably for intuitive and interactive Chinese learning. The web portal – MIMAS (http://scholar-milk.i2r.a-star.edu.sg/mimas3-v6/index.php/home2) is a fast tool for creating various sub-apps within the container app, simply by uploading and linking prepared contents. Teachers or educational content providers can then create smartphone apps easily for their own teaching topics, even without any programming skills. Different from Optical Character Recognition (OCR) technology, which is sensitive to the orientation of the character image and its background, our approach can recognize the character accurately from any angle. It does not matter whether the character is printed in uncommon fonts, on cluttered background or marred by drawings around it. The character recognition is reliable and fast without requiring the user to capture the image from the character at its upright position.

**Keywords:** Chinese learning, Image recognition, Augmented reality, Optical Character Recognition (OCR), App creation

## Background

The Chinese language is the most widely spoken language in the world. In recent years, Chinese language learning has been one of the fastest growing fields in second language education in the world. However, because of the complexity, variety and pictographic nature of its written characters, Chinese has been one of the most difficult languages to learn and remember especially when the characters have complicated strokes. Looking up the Chinese dictionary requires a few steps, including strokes counting, radical lookup and final character lookup. This can be tedious and time consuming for learners who do not know its pronunciation or pin yin representation. With the rise in the popularity of smartphones in recent years, many apps have been developed to enhance Chinese learning. Some of them work like a digital dictionary, requiring the users to input the characters in order to retrieve their meanings. To alleviate the cumbersome input process, advanced

technology like OCR is used to recognize characters from images captured by camera. However, OCR does not work well when the characters are printed on cluttered background, in uncommon fonts or in the presence of other drawings. In this paper we propose an alternative approach to overcome the above problems. With our approach, smartphone users can point their phone camera to the character from any direction. Regardless of the background and font of the character, it can be recognized quickly.

## Review of similar technology

Due to the widespread use of smartphones, the printed Chinese dictionary has been replaced by smartphone apps. Dictionary inputs are varied among these apps, with apps accepting inputs using pin yin and radicals, and even handwritings. Pin yin input has its limitation as most people may not know how to pronounce the character. That is why they look up the dictionary in the first place. Radical and handwriting input is tedious and time consuming. To solve this problem, OCR technology is

* Correspondence: yqli@i2r.a-star.edu.sg
Institute for Infocomm Research, A*STAR, Singapore, Singapore

Li *et al. Scientific Phone Apps and Mobile Devices* (2016) 2:4

Page 2 of 12

used to recognize characters on printed media to assist the learning of new characters. There are apps available in the App Stores using OCR technology, e.g., Pleco (http://www.pleco.com/pleco2ip.html), China Goggle (https://itunes.apple.com/us/app/china-goggles/id401162596?mt=8), Google Translate (https://itunes.apple.com/en/app/google-translate/id414706506?mt=8) and Hanping (https://play.google.com/store/apps/details?id=com.embermitre.hanping.app.reader.pro) etc. Figure 1 shows a screenshot of a result from Pleco app. The Pleco app requires the characters to be carefully positioned inside a box to determine the boundary for character segmentation. If the characters are fitted nicely inside the box and the background is clean, the recognition is robust. This requires some amount of coordination on the part of the users, who may not always be in an environment conducive for steady positioning of the characters (e.g. when the user is in a bus). Figure 2 shows the recognition results from China Goggle. A character is recognized correctly when it is placed upright inside the box. The recognition results are wrong when the character is not upright or if the background is not clean. Figure 3 shows how Google Translate works. Users first align the text inside the rectangle box before clicking on the camera button to capture the image. Users are then required to swipe across a region to highlight the text for recognition. As with other apps, recognition is affected when the background of the text is not clean. In this particular case, the character '酒' is not recognized as the pinyin representation is printed above it, as shown in the third screen capture in Fig. 3.

In summary, current technologies used to recognize Chinese characters on printed media are mostly based

on OCR. The OCR based apps for Chinese learning are reliant on the phone users to place the characters inside a box and at their upright position in order to get a good segmentation of the characters. This apps require some conscientious effort from the users who have to hold the phone steady and keep it at a specific position so that the targeted characters are placed nicely inside the box, while ensuring other characters or noisy background are excluded.

## Comparison between OCR and our approach for Chinese character recognition

From a technology point of view, OCR is a mature technology for recognizing characters from printed media with 99 % accuracy on Chinese newspaper (http://www.dlib.org/dlib/march09/holley/03holley.html), provided the characters are correctly segmented. However, the high recognition rate of OCR is largely dependent on the performance of its pre-processing steps. A typical mobile OCR engine goes through the following processes (http://www.abbyy.com/mobileocr/OCR_stages).

1. Image import and processing
   Image binarization is conducted to separate the text from the background. In this process, the skewed image must be corrected and the character orientation must be obtained.
2. Document analysis
   Characters are detected and segmented, then joined into word, lines of text etc.
3. Optical Character Recognition (OCR)
   Segmented characters are recognized using the special language and pattern definition. Image must be in good quality for fast processing and high accuracy. Sample characters in various fonts may have to be collected to train the OCR engine to recognize the same characters in different fonts.
4. Result processing
   The results can be further processed and verified depending on different applications. Lexicon can be used to verify the recognition result.

Hence, even before the recognition step, a series of pre-processing operations, including denoising, binarization, skew adjusting, page segmentation and character segmentation have to be conducted. The performance of the character segmentation is in turn affected by the complexity of the background of the characters and the acquisition of the image (Qiang et al. 2001; Yuan et al. 2013). If the characters are printed on colorful or noisy background, the character segmentation results may severely deteriorate. Images captured from a mobile phone camera pose more challenges than those acquired from a scanner as the characters captured by phone cameras


**Fig. 1** OCR in Pleco app

**Fig. 2** China Goggle: A character is recognized correctly only when it is placed upright inside the box

may undergo additional affine transformation and uneven lighting condition, making binarization difficult.

In summary, character segmentation is a critical process as the OCR engine requires a segmented clean and upright binarized image of the character as its input. When there is a row of characters or a paragraph of article on the image, segmentation can make use of the contextual information. On the other hand, if there are only one or two characters on the image, the orientation of the character is difficult to obtain. For a paragraph of article, lexicon verification is used to improve the recognition accuracy. Therefore, OCR is good for recognizing a paragraph or a page of article. It can achieve high recognition accuracy on documents printed on a media with uniform background.

The difference between our approach and traditional OCR approach for Chinese character recognition is that we use local descriptors to represent the Chinese text image. This local descriptor is invariant to image orientation, scale and partial occlusion. This method is normally used for textured rigid image recognition. From our experiments, we found that it also works on Chinese characters with certain complexity under some modification of the algorithms. This approach does not require

character segmentation and any other preprocessing of the image as what traditional OCR requires.

In our approach, we treat a Chinese character, word or phrase as an image pattern. We first detect salient keypoints on the whole image and then extract local feature at those salient points. This local feature is invariant to the image orientation, scale, partial occlusion and affine transformation. The local descriptors from the query image are matched against those in the reference images. The one with the most number of matches emerges as the results. As long as there are enough salient points on the image, this approach works well. If the number of salient points is not enough, as with simple characters such as "一", the image may not be recognized. The details of our method will be described in section Image Recognition Engine.

As we do not have to perform character segmentation as traditional OCR does, our approach can avoid errors caused by character segmentation. Therefore, our approach does not require the character to be upright and clean. It also does not require the character to be printed on a uniform background. With our method, Chinese text can now be printed on colorful background so that it looks more attractive to young students and children. More importantly, the characters can be captured from
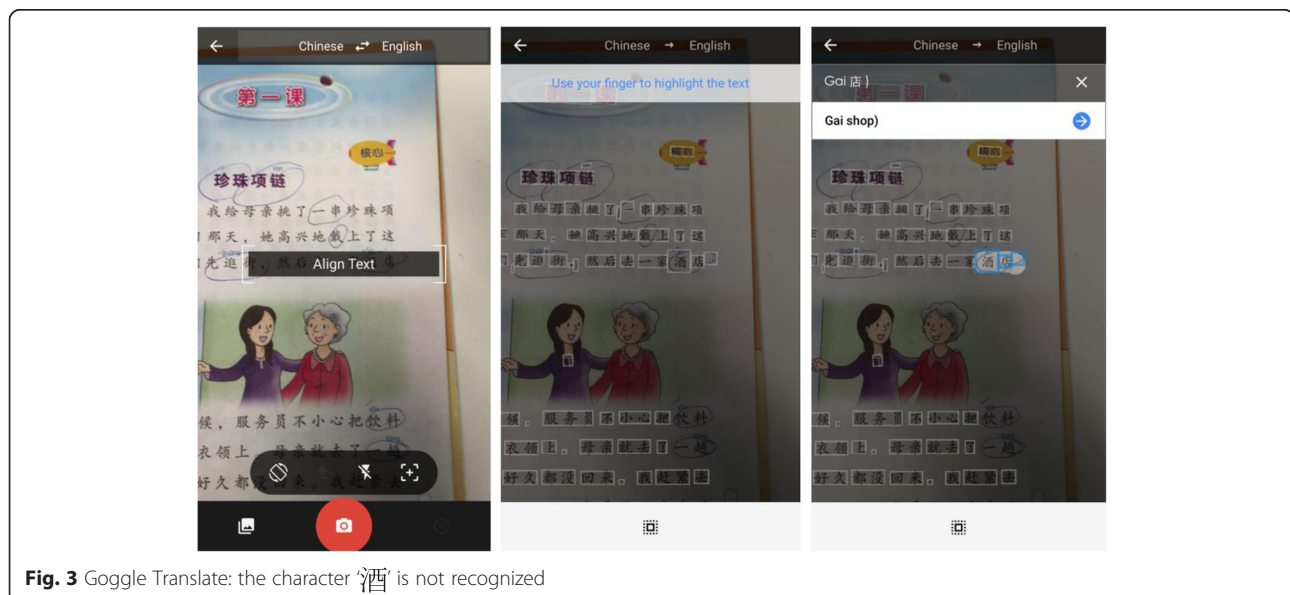


**Fig. 3** Goggle Translate: the character '酒' is not recognized

Li *et al. Scientific Phone Apps and Mobile Devices* (2016) 2:4

Page 4 of 12

any directions and angles. Mobile phone users do not have to conscientiously fit the characters into a box. Even though the surrounding characters are partially included or the targeted character is partially occluded, the recognition is still accurate. However, if the objective of the application is to recognize a paragraph of characters printed on newspaper or books etc., OCR technology is still the best choice. This is because the OCR engine is trained using large number of character samples in most common used fonts. Rather, our approach benefits applications that aim to recognize specific instances of a set of characters. Only one character sample per instance is needed for training, making the training process much easier and simpler. Even teachers without expertise on OCR can collect the character samples and upload them into our MIMAS platform to create a smartphone app for their students.

To summarize, our approach is superior to traditional OCR approach in the following situations:

1. When the Chinese characters are printed on colorful or non-homogeneous background.
2. When the Chinese characters are randomly placed, isolated, printed using uncommon fonts or partially occluded.
3. When an application (smartphone app) requires fast and accurate response because the targeted users do not have the capacity, time and patience to hold and maintain the phone at a specific position.
4. When the application (smartphone app) only has to recognize instances of a limited number of characters.

### Comparison between existing applications and our platform for Chinese learning

To the smartphone users, accuracy, usability and responsivity are important factors. Using our recognition approach, users do not have to spend time and effort to place the characters upright inside a designated box. They can point the phone camera at the characters in any orientation within a range of distance, as opposed to fixed orientation and distance required by some apps.
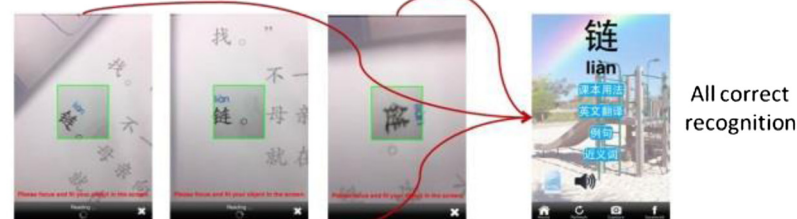
Therefore, it would be easier for the users to get a quick response. Figure 4 shows that our app can recognize a Chinese character in various directions. It also works on cluttered backgrounds such as with pin yin on top of the character.

Apps can be created from our web portal, which configures results of the dictionary entry. Results can be displayed almost immediately when the users point their phone camera at the character, hence saving time for the learners. Traditional tools often produce wordy results and may not appeal to some readers. Our app then allows contents in multimedia forms such as web pages, audio or video clips or 3D graphics to explain the meaning of the characters. Displays with augmented reality are provided in our system. The display of the results is customizable based on the instructional design of the teachers or content providers. The web portal is easy to use and apps can be created quickly and automatically without programming, saving time and cost on app creation.
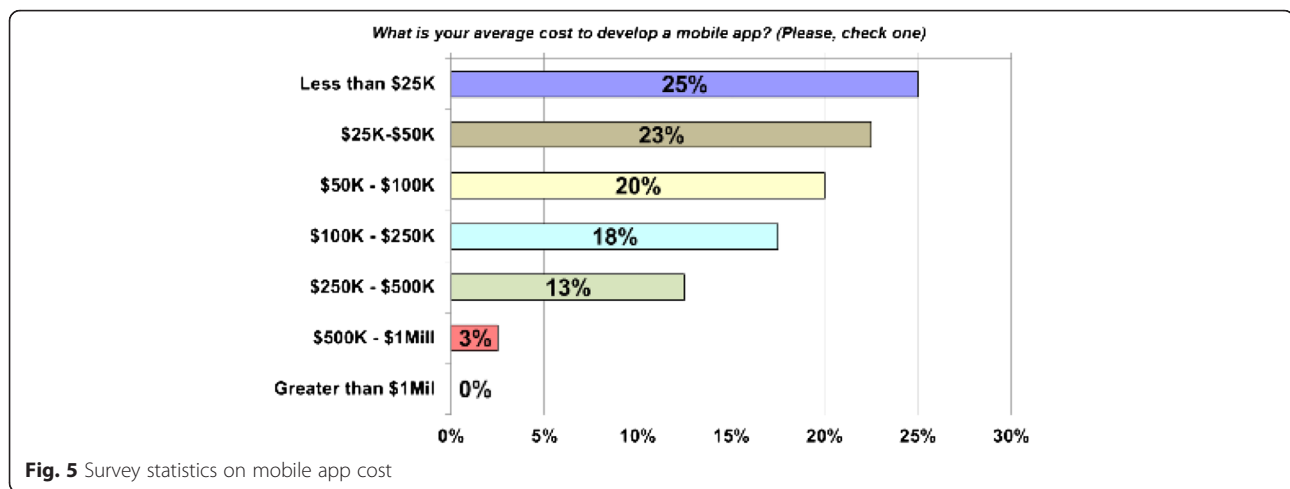
### Implementation

Creating a smartphone app is time consuming and requires mobile programming skill. It is expensive and time consuming for content providers such as teachers to engage companies to create the app. Apps with advanced augmented reality technology are more complicated and require extensive expertise such as image recognition, camera motion tracking, 3D graphic rendering etc. According to a survey statistics from 2014 MGI Research (http://www.mgire-search.com/Mobile-Computing/how-much-does-a-mobile-app-cost-how-long-does-it-take.html), 43 % of the mobile apps cost \$25–\$100 K (Fig. 5), and 96 % of the apps require more than one month to develop (Fig. 6). Therefore, a platform that simplifies the smartphone app creation process will be able to help content providers to save cost, time and efforts.

We have developed a platform with a simple Graphic User Interface (GUI) for the users to configure or customize the mobile app that they intended to create. Teachers can make use of the platform to create mobile apps easily for students to learn Chinese from printed



Correct recognition in any cases either the character is rotated or the background is cluttered.

**Fig. 4** Recognition results from our app

Li *et al. Scientific Phone Apps and Mobile Devices* (2016) 2:4

Page 5 of 12



**Fig. 5** Survey statistics on mobile app cost

media. Users only need to upload one sample image for each Chinese character and associate multimedia contents to each of the image. A mobile app can be created without any programming. The high level system architecture is shown in Fig. 7.
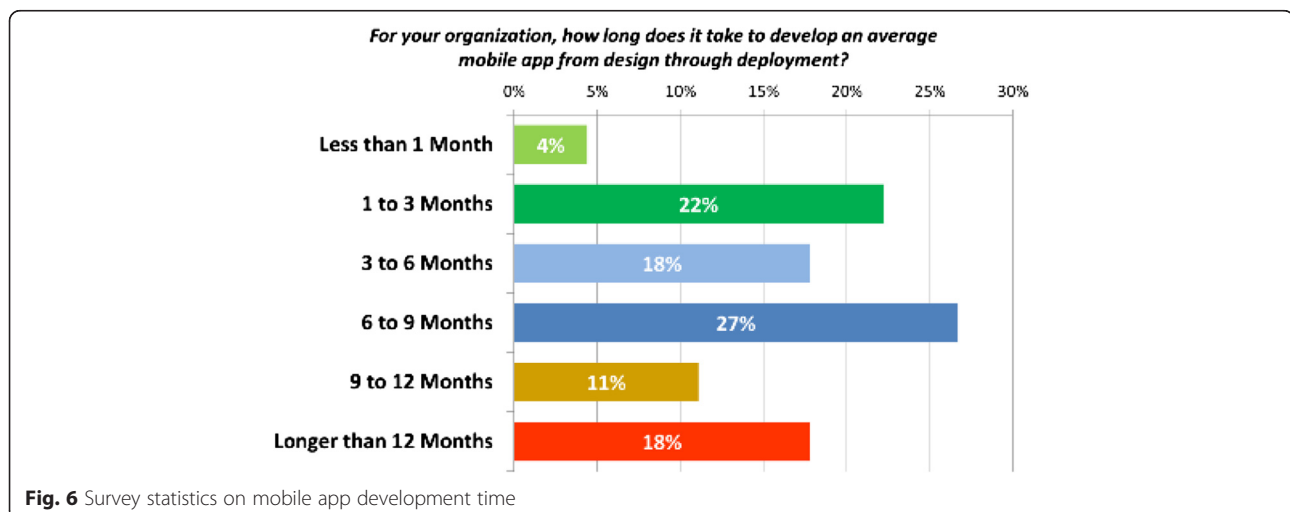
The app creation platform is able to create smartphone apps quickly and easily with prepared contents. Users only need to upload one image for each character, word or phrase in a zip file and associate each character image with its corresponding contents to explain its meaning. The training of the recognition engine is then executed in the background. After it is completed, the app will be created automatically. Using the created app, the users can point their phone camera at a Chinese character, word or phrase, the corresponding multimedia content will be displayed immediately, explaining the meaning of the queries. The multimedia contents can be a webpage, image, video, audio, 3D graphics or any of these combination. With our camera motion tracking

and pose estimation software, advanced Augmented Reality (AR) apps can be created to blend 3D graphics or videos into the physical objects with motion tracking capability. The created apps can be run either offline without internet connection or online as a client, connecting to a server to access up-to-date live contents.

### User Interface design

The GUI of web portal is designed to allow users to login, upload target images, associate multimedia contents to the target images etc. Figure 8 shows the user interface after user login. At this stage, users can choose to read the user guide, create a new app or review an existing app.

To create a new app, users have to upload the target images captured from one or more Objects of Interest. To speed up the image uploading process, the system allows users to zip the folders which store images captured from various Object of Interest and upload all images once.



**Fig. 6** Survey statistics on mobile app development time

Li *et al. Scientific Phone Apps and Mobile Devices* (2016) 2:4

Page 6 of 12



**Fig. 7** High level system architecture

This saves uploading time as opposed to uploading the images one by one. More Objects of Interest can be added in later by uploading another zip file. Figure 9 shows the GUI after a zip file of the images is uploaded.

For each Object of Interest, a thumbnail image is displayed. A set of fields is provided for the user to input the content links or upload contents associating with the highlighted Object of Interest (Fig. 10).



**Fig. 8** GUI after user login the system

Li *et al. Scientific Phone Apps and Mobile Devices* (2016) 2:4

Page 7 of 12



**Fig. 9** GUI after users upload target images (Object of Interest)

For 3D contents, the user can configure the position and pose of the 3D model in relation to the physical object (Fig. 11) after uploading.

To enable better user experience and flexibility of various use cases, there is a "Settings" tab in the web portal (Fig. 11). Under the 'Settings' tab, there is an "IR Configuration" option, allowing the user to configure whether the image recognition will be performed on the mobile device or on the remote server. If the user chooses "Offline Mode", the image recognition will be performed on the mobile device without sending image to a remote server. Some of the contents will be downloaded to the device and then the app can be executed without Internet connection. If "Online Mode" is chosen, the image recognition will be performed on a remote server. In this case Internet connection is required for the app to send image to the server for recognition.

Under the "Camera Configuration" panel, users can configure different camera modes to capture a single image or continuous video stream for recognition and activation of the contents. Each mode has its advantages and disadvantages, depending on the use case scenarios and the presence and cost of Internet 3G or WiFi connection. Under "Normal" mode, one image will be captured for image recognition once at the click of the camera button. If the image recognition is conducted at a remote server, the app only sends one image to the server for recognition. It will minimize the cost of sending image data to the server. In "Auto" mode, when the user clicks on the camera button, images will be captured continuously until an image is recognized and the content is displayed. With this configuration, the user does not have to touch the camera button multiple times if the image is not recognized in the previous captures. For the "Live" mode, the camera will keep capturing



**Fig. 10** Fields for content association

Li *et al. Scientific Phone Apps and Mobile Devices* (2016) 2:4
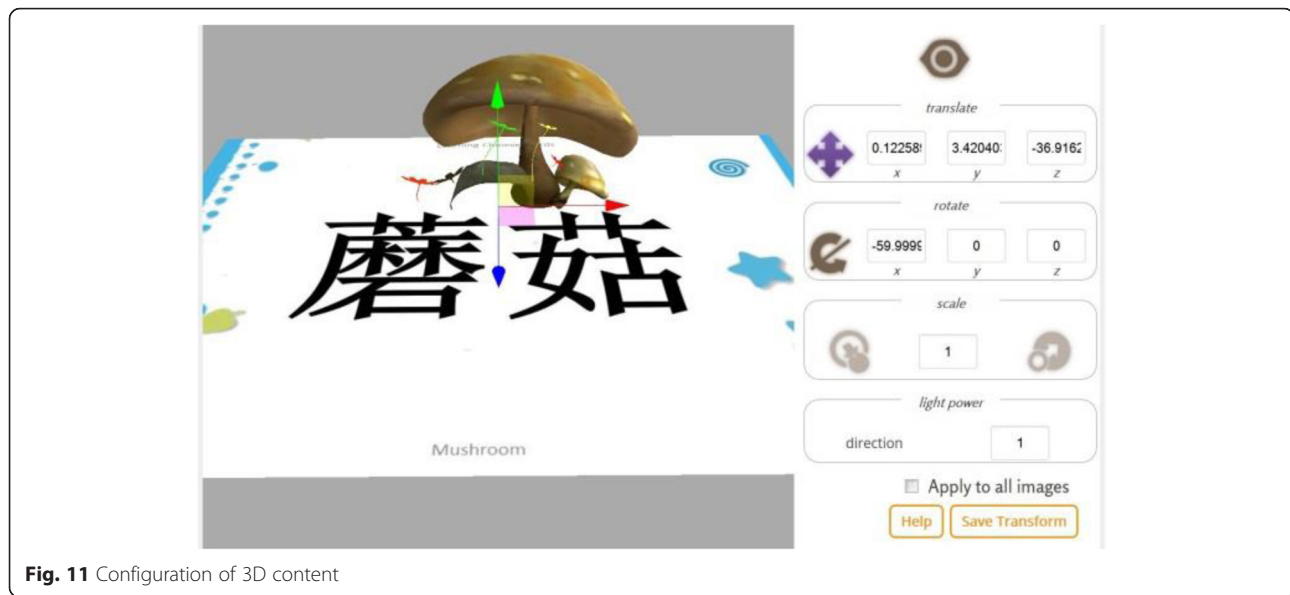
Page 8 of 12


**Fig. 11** Configuration of 3D content

images for recognition even though an image has been recognized and the content is displayed. The display keeps refreshing when the camera moves from one object to another, which can be used for quick identification of various objects appearing in the camera view. If the image recognition is performed at a server for the "Auto" and "Live" mode, the Internet 3G or WiFi data usage will be intensive as the app is continuously sending data to the server.

By default, the image recognition engine will use the full screen image as input for recognition. If the Object of Interest is far away and thus becomes too small in the camera view, a full screen image may include a lot of background which is not relevant to the Object of Interest. In this case image cropping can be enabled, displaying a green rectangle in which the app user can place the Object of Interest.

After all configurations and settings are done, clicking on the "Finalize App" (Fig. 12) will activate a series of automatic processes to create the app. The processes include,

- Checking for image quality and validity among the user-uploaded images as users may not be aware the suitability of an image for image recognition. This automatic process will inspect each image to check whether it is blurred, texture-less or a duplicate of another category. If the image is unsuitable or invalid, the users will be informed.
- Performing feature extraction and model creation for the uploaded images.
- Setting up the user interface for the client app, e.g., the camera mode and the image recognition mode. This enables a dynamic and flexible mobile app with

customized user interface in order to accommodate the preferences of various users.

## Image recognition engine

In the mobile app, the content to be displayed is driven by the mobile device's camera, where its captured images are fed to and identified by our image recognition engine. To facilitate various user requirements, we ported our image recognition software to different platforms including Linux, Windows, iOS and Android. Therefore the image recognition engine can run on either smartphones or server computers. If it runs on the smartphone, the app does not need to send the query image to the server. The image recognition will be performed on the smartphone. If it runs on the server, the app has to send the query image to the server where image recognition is performed. In both cases, the app will retrieve and display the corresponding contents after obtaining the recognition result. As mentioned in 2.1, the image recognition engine can be configured to run on the smartphone or the server, depending on whether Internet connection is available, and the number of images to be recognized. If Internet connection is not available, the image recognition will be performed on the smartphone. If the number of images to be recognized is large, memory availability on the smartphone may be an issue and hence it is recommended that the image recognition be performed on a remote server.

Before the image query, a set of sample images has to be collected. Features are then extracted and stored in a database. During the image query, the system extracts features from the query image and then matches them with the ones in the database. Geometric validation is performed to verify the correctness of the matching (Fig. 13).
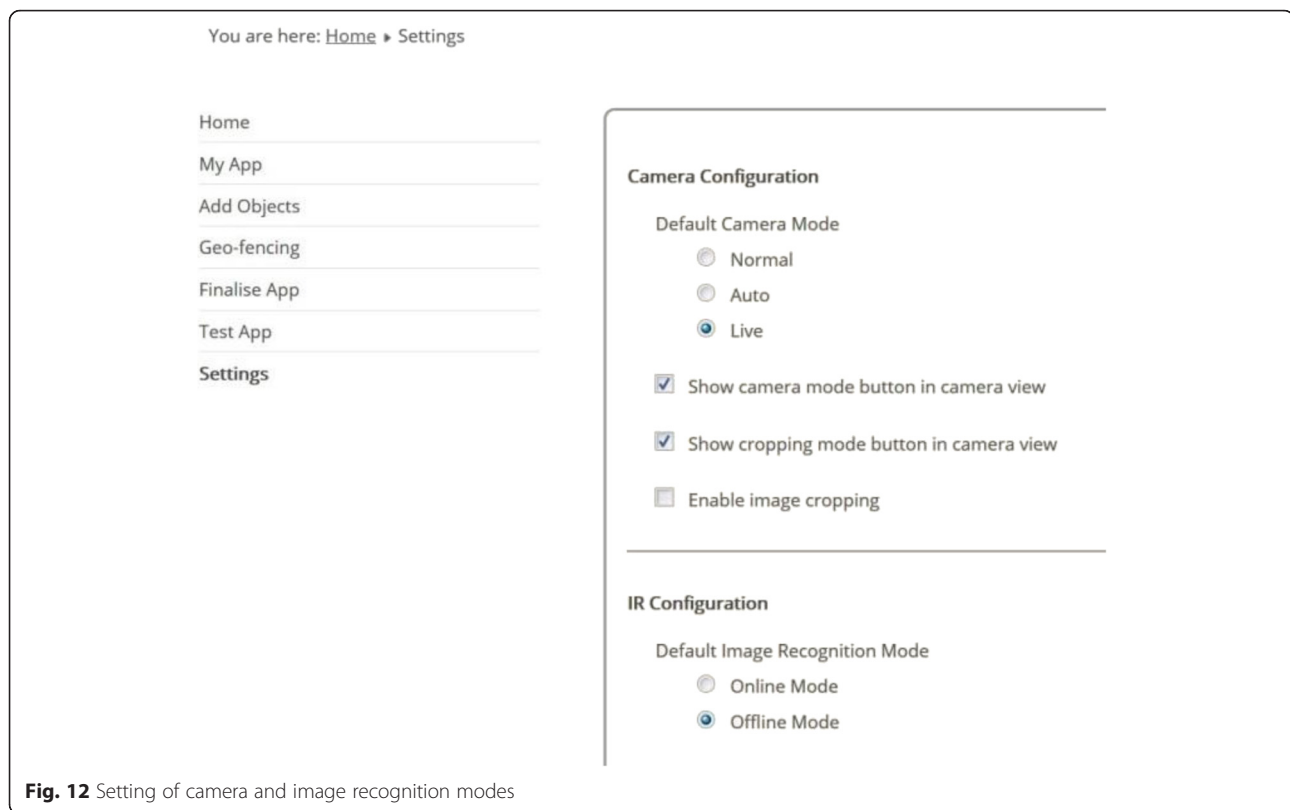
Li *et al. Scientific Phone Apps and Mobile Devices* (2016) 2:4

Page 9 of 12



**Fig. 12** Setting of camera and image recognition modes

The top matching image will be derived according to the number of matching feature keypoints (Fig. 14). The matching points are further verified based on their spatial relationship using the Longest Increasing Subsequence (LIS) (Gede et al. 2012). Due to its simple calculation, LIS for geometric validation consumes less memory and its computation is faster than the popular RANSAC homography approach. It is suitable for on device image matching so that markerless mobile augmented reality with real time camera motion tracking can be realized. The principle for the LIS method is that a set of matched keypoint pairs are geometrically consistent if their order is the same in both the query image and the database image. This is because the geometric order of the keypoints does not change when we translate or scale the images. Therefore, we can determine whether a query image matches a database image based on the largest subset of the matching keypoints by calculating the longest increasing subsequence on the spatial coordinates of the keypoints. In this way the outliners can be removed from the true matching



**Fig. 13** Image recognition engine

Li et al. Scientific Phone Apps and Mobile Devices (2016) 2:4

Page 10 of 12



**Fig. 14** Corresponding matching feature keypoints

keypoints. The image matching result will be more reliable.

We can go further with the image recognition results by performing camera motion tracking and pose estimation to obtain accurate position and orientation of the camera with respect to the Object of Interest. Our system uses a hybrid feature and template based approach to perform the tracking (Gede et al. 2015). Tracking allows us to create augmented reality scenes by rendering a virtual object positioned at the estimated pose against the background of camera image stream, giving the impression that the virtual object is part of the real environment. With camera motion tracking, we can then expand our range of contents to 3D graphics and video overlays, making our Chinese learning app more interesting.

## Results

We created an app using our platform and image recognition engine for 300 new Chinese characters from the primary 6 textbook used in Singapore schools. We first captured the character images from the textbook and uploaded them into the platform. A webpage explaining the meaning of each character and an audio clip of its pronunciation are then associated with the character image. After these are done, with one button click, a Chinese character learning app was created. Figure 15 shows a few screen captures of the app.

We scanned the 300 new Chinese characters from the primary 6 textbook using an iPhone 5 with the MimasAR app installed. For each character, the correct dictionary entry was retrieved, meaning virtually 100 % correct



**Fig. 15** Screen captures of the Chinese character learning app

Li *et al. Scientific Phone Apps and Mobile Devices* (2016) 2:4

Page 11 of 12

recognition rate for these 300 characters. We have recorded a video and uploaded it to YouTube (https://www.youtube.com/watch?v=RZ9_1C85Gh0).

## Discussion

Above testing was conducted in an office environment with good lighting condition. If the testing is conducted under dim lighting condition, the accuracy may be affected. Furthermore, if we have more characters in the database to be recognized, the accuracy may be lower as well. We will test more characters in the future to understand how the database size will affect the accuracy.

## Conclusions

We developed a web portal and a container app with image recognition and augmented reality technologies for quick and easy Chinese character learning. When the app users point their phone camera at a Chinese character, the character is identified immediately and the corresponding multimedia content explaining its meaning will be displayed. The multimedia contents can be webpage, image, video, audio, 3D graphics or any of these combinations. This technology saves the learner's time on dictionary lookup, and makes Chinese learning intuitive and interesting.

Different from the OCR technology, our approach can recognize Chinese characters accurately and quickly from any angle. It does not matter whether the character is printed on cluttered background, or marred by drawings around it. Unlike OCR, our approach only requires one image sample of the character to train the recognition engine.

Our current system allows people to create smartphone apps within an existing container app. In future, we plan to further develop the system to create self-branding apps, which are app packages with customized user interface ready to be submitted to app stores.

## Availability and requirements

Our automatic app creation platform is available at:

http://scholar-milk.i2r.a-star.edu.sg/mimas3-v6/index.php/home2

Evaluators can register for an account using an email address to test or evaluate the platform free of charge for one month. The platform is available for commercial license if anyone is interested in using it for their business. In fact, it is licensed to a few companies.

Our iOS and Android container app are available for free download at Apple App Store and Google Play Store below.

https://itunes.apple.com/sg/app/mimasar/id988692600?mt=8

https://play.google.com/store/apps/details?id=sg.edu.astar.i2r.mimasar&hl=en

- Project name: Mobile Interactive Multimedia Access System (MIMAS)
- Project home page: http://scholar-milk.i2r.a-star.edu.sg/mimas3-v6/index.php/home2
- Operation system(s): Linux or Windows for the MIMAS platform. iOS and Android for the iOS and Android container apps
- Programming language: PHP, Java Script, Objective-C, Java, C/C++
- Other requirements: users have to agree on a Trial Web Service Agreement at http://scholar-milk.i2r.a-star.edu.sg/mimas3-v6/index.php/component/users/?view=registration in order to use the platform free of charge for 30 days.
- License: commercial license with license fee
- Restrictions to use by non-academics: either academics or non-academics users can register for an account to use the platform free of charge for 30 days.

## References

Gede Putra Kusuma N, Fong Wee T, Li Y. Hybrid feature and template based tracking for augmented reality application, proceeding of Asian Conference on Computer Vision (ACCV) workshop, vol. 9010. 2015. p. 381–95.
Gede Putra Kusuma, Attila Szabo, Yiqun Li, and Jimmy A. Lee: Appearance-Based Object Recognition Using Weighted Longest Increasing Subsequence. In

Li *et al. Scientific Phone Apps and Mobile Devices* (2016) 2:4

Page 12 of 12

Proceedings of the 21st International Conference on Pattern Recognition (ICPR), p. 3668–3671, Tsukuba Science City, Japan, November 2012.

Qiang H, Yong G, Zhi-Dan F. High performance Chinese OCR based on Gabor features, discriminative feature extraction and model training. Proceedings of the Acoustics, Speech, and Signal Processing, vol. 3. 2001. p. 1517–20 (ICASSP '01).

Yuan Mei, Xinhui Wang and Jin Wang, A Chinese Character Segmentation Algorithm for Complicated Printed Documents, International Journal of Signal Processing, Image Processing and Pattern Recognition, Vol. 6, No. 3, 2013.