



**User Guide for
Optical Character Recognition System
for Sinhala
(Alpha Release)**

1st October 2005

**PAN Localization Project
Language Technology Research Laboratory
University of Colombo School of Computing
Sri Lanka**

Table of Contents

1. Introduction-----	3
2. Running the Application-----	3
3. Using Application-----	5

1 Introduction

This document explains how to use the Optical Character Recognition (OCR) system developed to recognize Sinhala Characters. This OCR system requires totally noise removed and skew-corrected images in jpg format as its input. It writes its output to Unicode text files. The current application supports many widely used fonts. Users are expected to select the appropriate font before starting recognition. Before using this application please read the instruction given in the following sections and the Appendices carefully.

2 Running the Application

After the successful installation of Optical Character Recognition system for Sinhala (Refer installation guide for Optical Character Recognition system for Sinhala to install the application) and all dependencies, follow the steps below to start the application.

- I. Run the application from Windows start menu at location **Start -> All Programs -> Sinhala OCR (Alpha Release) -> Sinhala OCR** illustrated in Figure 2.1. If the link for “**Sinhala OCR (Alpha Release)**” is missing in the start menu please reinstall the application.



Figure 2.1: Running the “Optical Character Recognition System for Sinhala” from Windows start menu

II. As the application starts, it shows the main window as shown in Figure 2.2.



Figure 2.2: Optical Character Recognition system for Sinhala, at start up

3 Using Application

3.1 Open the source image

In order to open the image or images, follow the steps given below:

- I. Click the “Open” button. Then the open window will appear as shown in Figure 3.1.1

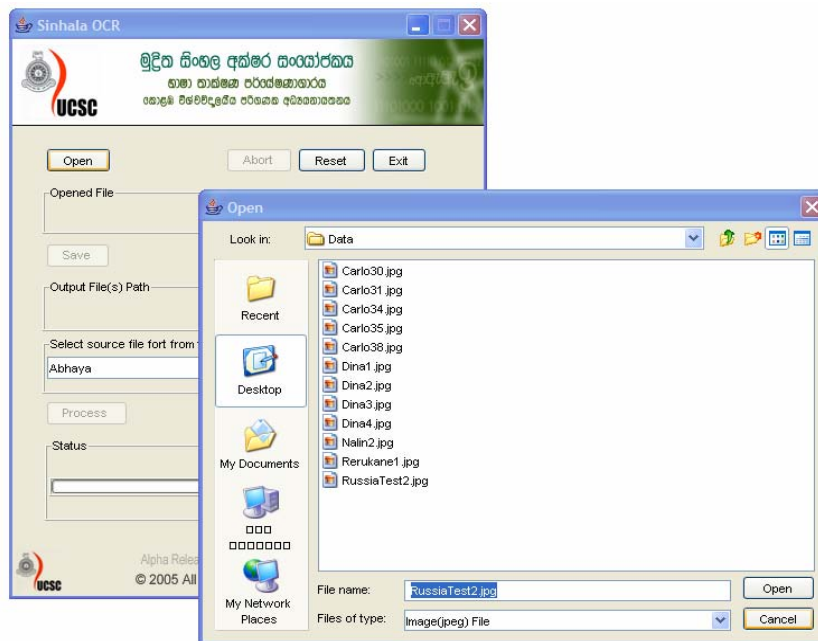


Figure 3.1.1: open the source image(s)

II. As shown in Figures 3.1.2(a) and 3.1.2(b), choose one or more images for processing

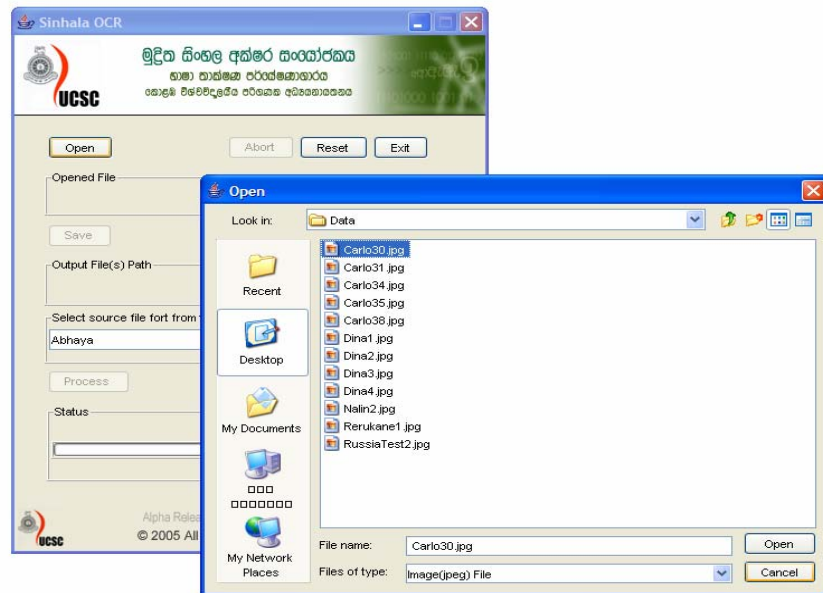


Figure 3.1.2 (a): Choose an image file

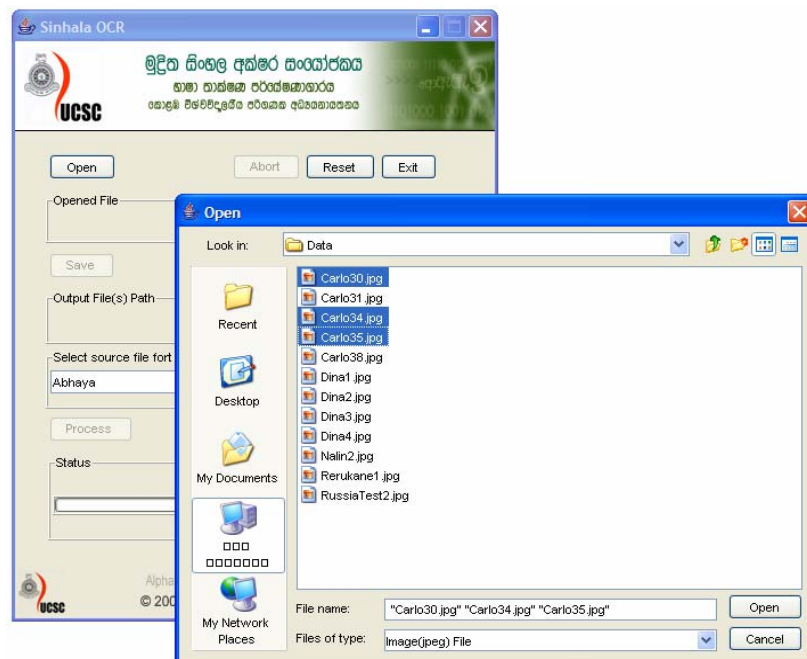


Figure 3.1.2 (b): Choose more than one image file

- III. After an image (or images) has (have) been selected the application shows the path of the selected image and the default path where the output will be saved. See Figure 3.1.3

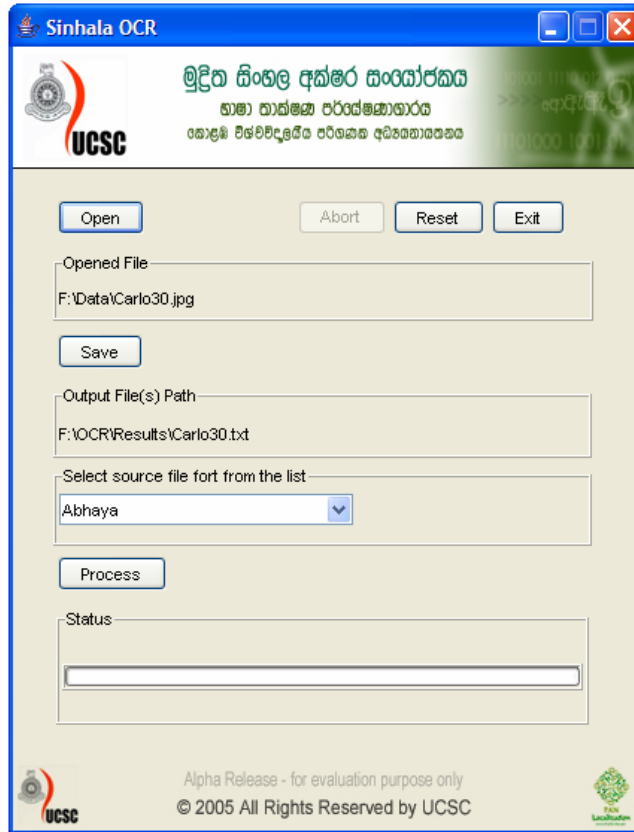


Figure 3.1.3: Application shows the source and output file paths

- IV. If the user wants to save the output at a different location select the location by clicking Save button, see Figure 3.1.4

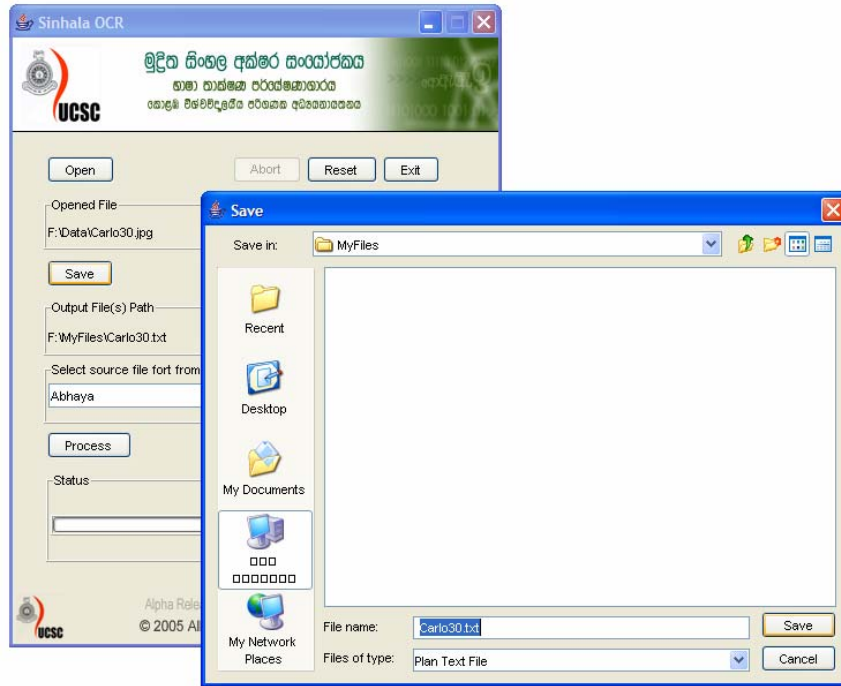


Figure 3.1.4: Output file location can be specified by the user

3.2 Select the font

In order recognize characters accurately the user has to specify the font in the source image. If the required font is not given in the list of fonts appear in the application select the closest one. See Figure 3.2.1. Samples are given in the Appendix 2.

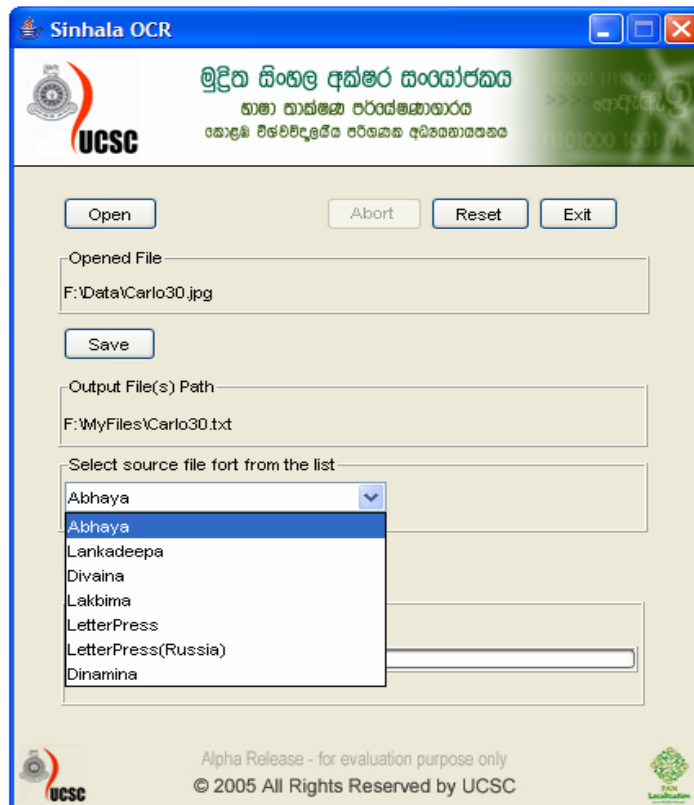


Figure 3.1.3: Select the appropriate font

3.3 Recognition

- I. Recognition of characters in a given image(s) will start as the Process button is clicked. See Figure 3.3.1.

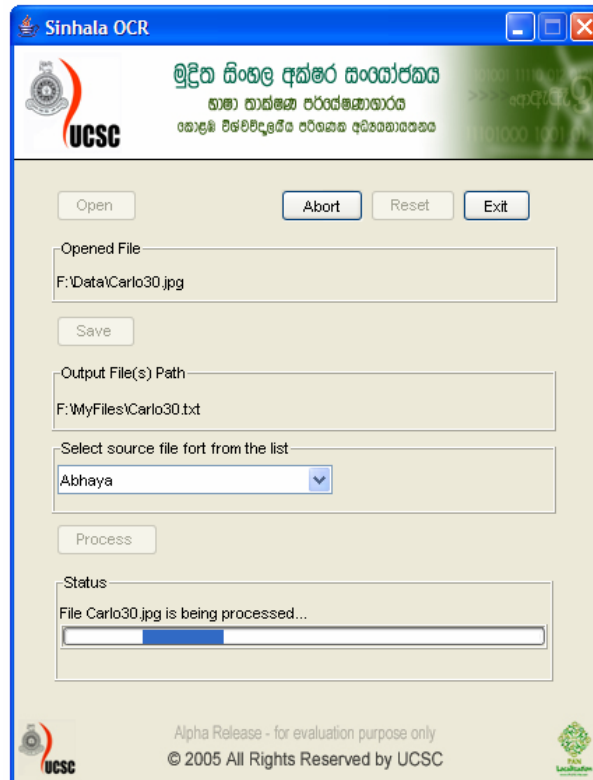


Figure 3.3.1: An image is being processed

- II. The application gives the message shown in the Figure 3.3.2 when the image is successfully recognized.



Figure 3.3.2: Image has been successfully recognized

- III. The application gives the message shown in the Figure 3.3.3 when the image is being processed has too much of noise. The recognized part of the image will be given as the output. User has to clean the image and retry.

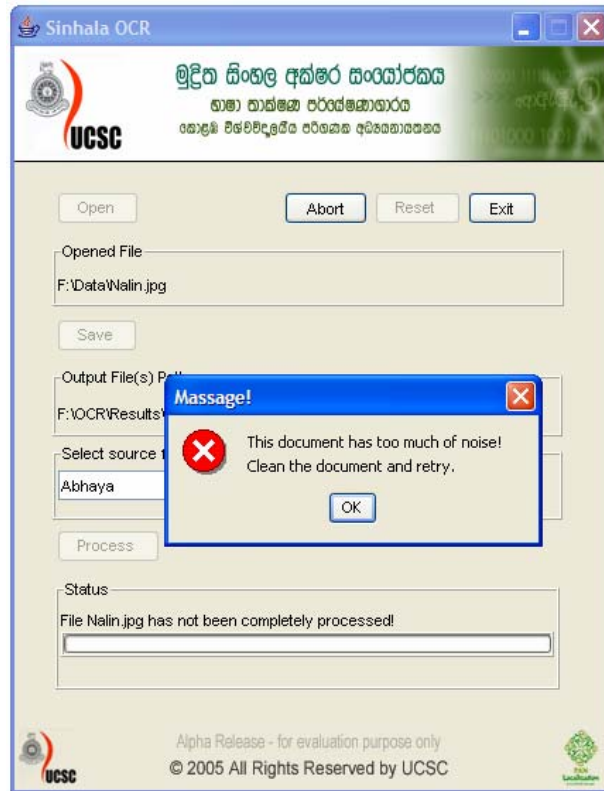


Figure 3.3.3: Image has too much of noise

- IV. Application tells the user the number of images successfully recognized out of given input files. See Figure 3.3.4.

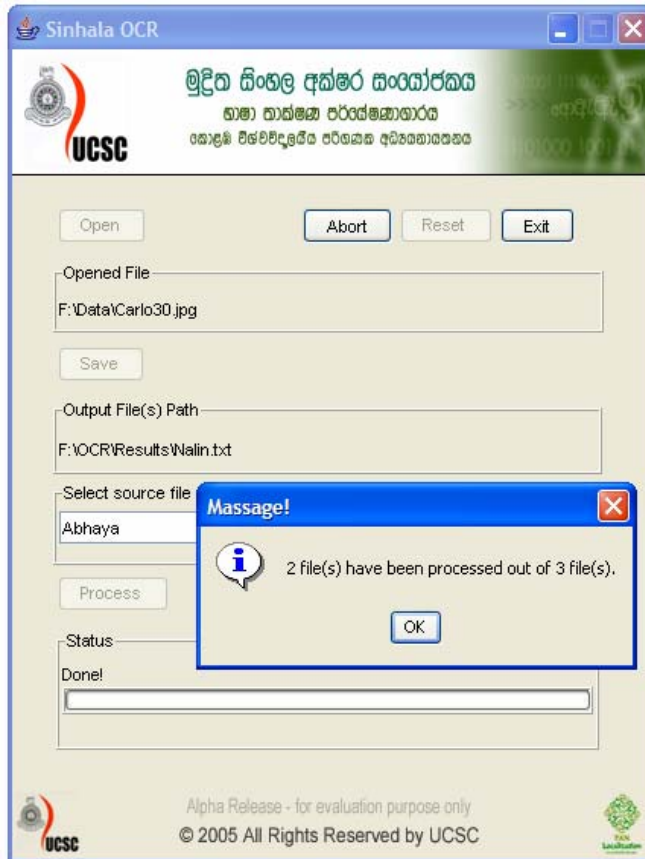


Figure 3.3.3: Number of successfully recognized files will be given by the application

3.4 Other

3.4.1 Reset

Click Reset to refresh the application, see Figures 3.4.1 (a) and (b):



Figure 3.4.1(a): Before resetting

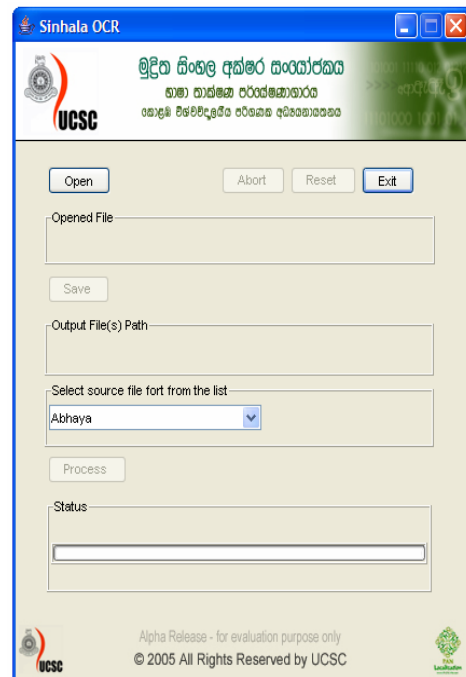


Figure 3.4.1(b): After resetting

3.4.2 Abort the application

User can stop the application at any stage by clicking Abort button. The consequent messages are given as shown in Figure 3.4.2. The user can decide whether to abort or not.

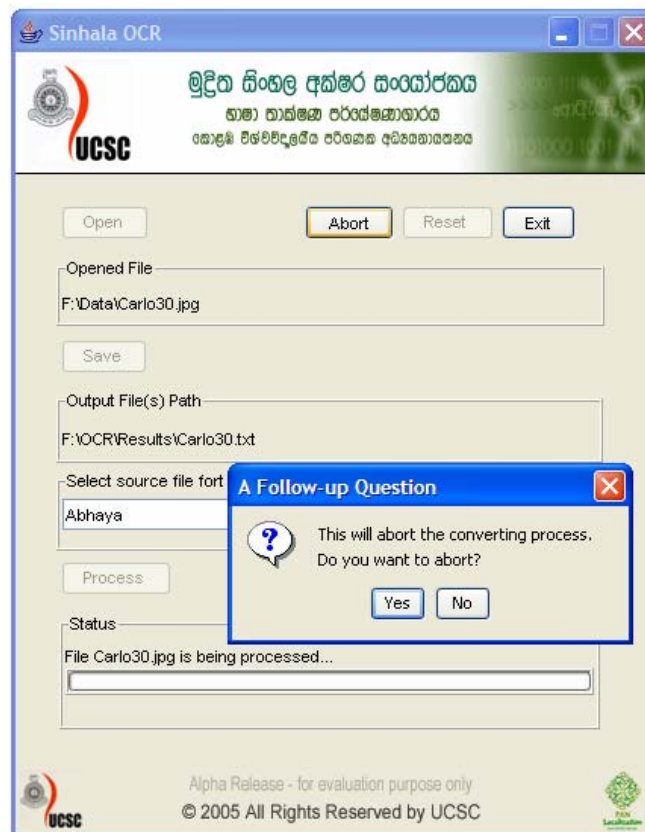


Figure 3.4.2(a): Before resetting

3.5 Close the application

Click Exit to close the application. See Figure 3.5.1.

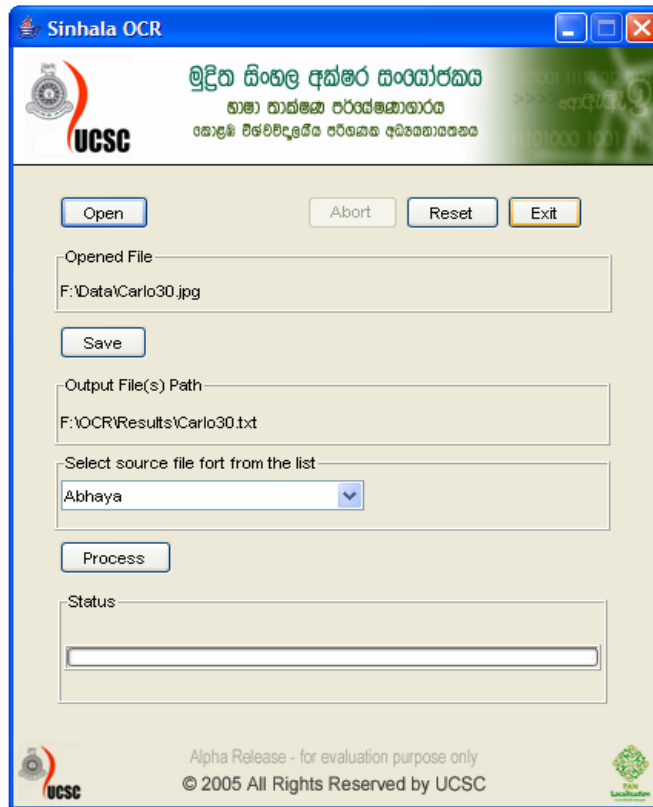


Figure 3.5.1: Click Exit to close the application

Appendix 1

Instructions to prepare input images

- Select a page from a book or a newspaper to convert into text.
- Use a standard scanner to obtain the image that is needed to be converted into text.
- Make sure to place on the scanner the paper in a manner it is not skewed. See Figures (a) and (b)

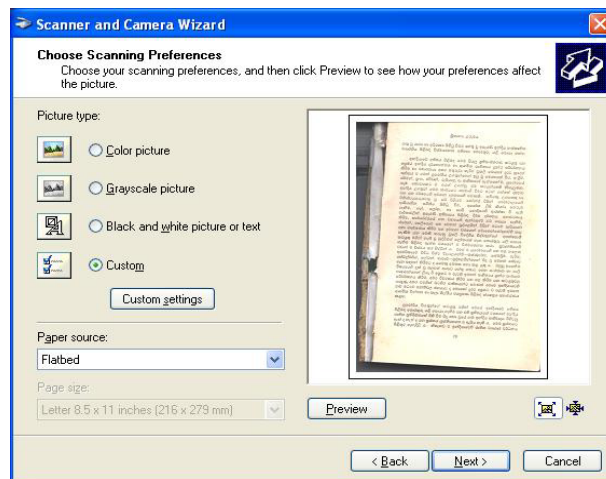


Figure (a): This image in too skewed

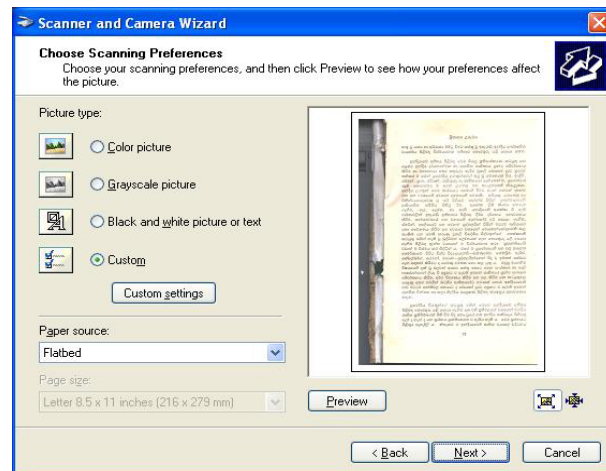


Figure (a): This image in not skewed

- Set the Resolution value to 600 dpi (dots per inch) in the scanner for better recognition.
- Obtain a preview of the image and select the relevant part of the image. See Figure (c)

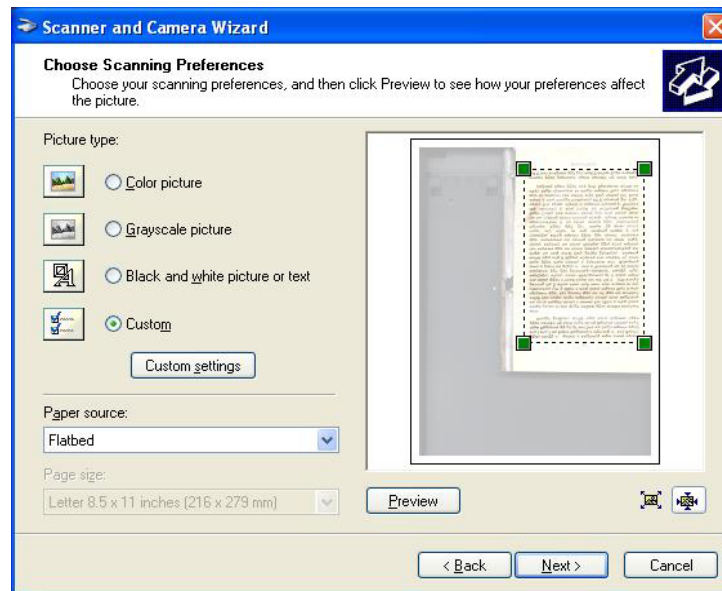


Figure (c)

- Save the image as a jpg image. (note: current OCR supports only jpg images)

- Remove noise in the image, if present, using a standard image processing application such as Adobe¹ Photoshop. See Figure (f)

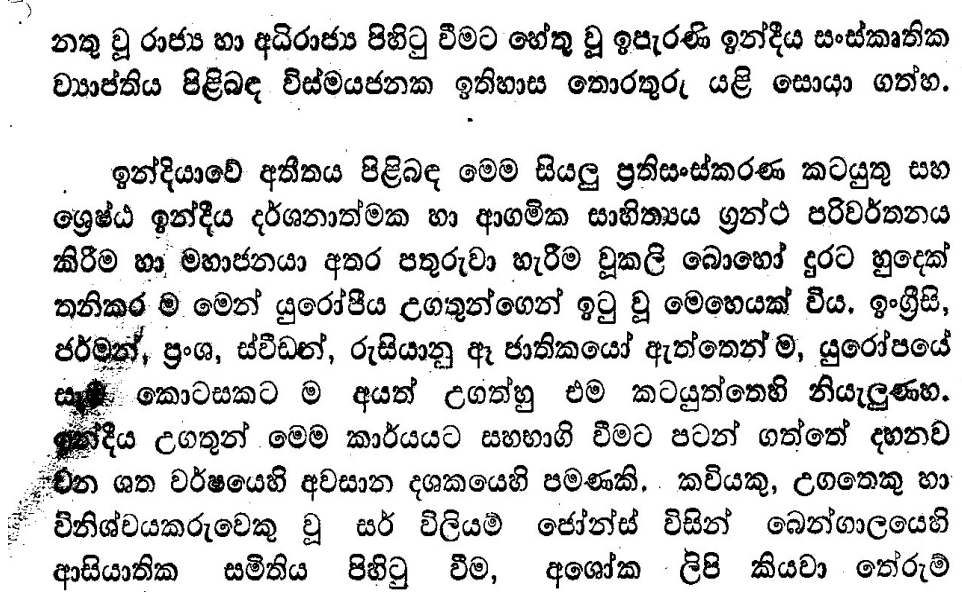


Figure (f)

¹ Adobe is a registered trademark of Adobes Systems Incorporated

Appendix 2

Supported fonts

The OCR supports widely used set of fonts. The available fonts are given in the application. User has to select a closest font from the list when a different font has been used in the document. Examples are given below for each font given in the application.

මේ ප්‍රශ්නය ගැන සිතන විට මගේ මනසේ මැවී පෙනෙන්නේ එක්තරා ජවනිකාවකි. ස්ථානය රෝහල් වාට්ටුවකි. එහි එක් ඇඳක වැතිර ඉවසා ගත නොහැකි දුක් වේදනාවෙන් මරලතෝනි නගන්නේ හැටපස් හැවිරිදි රෝගියෙකි. ඔහු පෙළනු ලබන්නේ මුත්‍රාශයේ හට ගෙන දසන පැතිරී ගිය දරුණු පිළිකාවකිනි. ඔහු අසාධ්‍ය තත්ත්වයක සිටින බව ඔහුගේ ආදරණීය බිරිඳ සහ වැඩිමහලු දරු තිදෙන දනිති. ඔහු ද ඒ වග දනී. ඔහුගේ වේදනාව නැසීම පිණිස වෛද්‍යවරුන් කරන ප්‍රතිකර්මවලින් ඔහුට සහනයක් නො ලැබේ. ඔහු ඒ දිනවල කළේ එක ම එක ආයාචනයකි. එය මෙසේ ය: “අනේ දොස්තර හාමුදුරුවනේ, මට මැරෙන්න බේත් ටිකක් විදින්න. තමුන්නාන්සේට බුදු බව අත් වෙයි.” මෙසේ ඔහු තොරතෝවියක් නැති ව මොර ගැවේ ය. ඔහුගේ බැගෑපත් ඉල්ලීම ඉටු වුණේ නැත. අපේ සම්මත නීතියට අනුව ඔහුගේ ඉල්ලීම ඉටු කිරීම බරපතල දඬුවම් ලැබිය යුතු අපරාධයක් බැවිනි. එපමණක් ද නො වේ. ඒ කාලයෙහි දී, එනම් වර්ෂ 1960 දී, සිය පණ නසා ගැනීමට යමකු තුළ ජනිත වන සිතිවිල්ල මා සැලකුවේ පාපකාරී එකක් ලෙස ය. එබැවින් පාප කර්මයක් කිරීමට ඔහුට රුකුල් දීමට මගේ හෘදය සාක්ෂිය මට අවසර නො දුන්නේ ය. තව දින කිහිපයක් ඔහු අතිශය පීඩාකාරී දුක් වේදනා විඳ මිය ගියේ ය. ඒ දුක විඳීමෙන් ඔහුට හෝ අන් කිසිවකුට හෝ කිසි දු යහපතක් සිදු නො වී ය. ඒ දින කිහිපය තිස්සේ ඉවසිය නොහැකි දුක ඔහු පෙළු සැටි මට කිසි කලෙක අමතක නො වේ. මා එවැනි තත්ත්වයකට අනාගතයෙහි දී පත් වුව හොත් මගේ පණ නසා ගැනීමට මම පසුබට නො වෙමි.

සිය දිවි හානි කරගැනීම අරඹයා බ්‍රිතාන්‍යයේ වෛද්‍යවරුන්

Figure (a): Abhaya

වෙහෙර විහාරවල් තැනින් තැන පැන නැගිණ. රජ මාලිගාවේ සංස්කෘතියත් බුද්ධාගමෙන් පෝෂණය ලැබීය. ජපන් විත්‍ර ශිල්පිහු චීන සාම්ප්‍රදායික වෙළුම් විත්‍ර කලාව අනුගමනය කරමින් බුද්ධ චරිතය නිරූපණය කළහ. අඩි පනස් තුනක් උස වූ නරා දැඩිබුන්සු නමින් හැඳින්වෙන ප්‍රසිද්ධ මහා බුදු පිළිමය ද නෙළන ලද්දේ මේ අවධියෙහි ය.

අනතුරුව පහළ වූ හෙයියන් හෙවත් කොන්තො යුගයෙහි (794-1185) බුද්ධාගම ජපානයෙහි ආදියෙහි පටන් පැවති ඕන්තෝ (දේව මාර්ගය) ආගම හා සම්මිශ්‍රණය වී වඩා දේශීය ස්වරූපයකින් වර්ධනය වන්නට විය. ජපන් බුද්ධාගමේ බලගතු සාධා දෙකක් වන තෙන්දයි නිකාය හා ඕංඕ නිකාය ද චීනයෙන් ගෙන එන ලදුව ජපානයෙහි පිහිටුවන ලද්දේ ද එම යුගයෙහි ය.

එහෙත් ජපන් සංස්කෘතියේ මූලස්තම්භයක් ලෙස සැලකිය හැකි ධ්‍යාන (චීන: චං, ජපන්: සෙන්) බුද්ධාගම එහි ගෙනෙන ලද්දේ කමකුරා යුගයෙහි (1185-1333) ය. ඒ යුගයෙහි ජපන් සමාජයෙහි ද වැදගත් විපර්යාසයක් සිදුවිය. එනම් සමුරයි නමින් හැඳින්වෙන ක්ෂත්‍රීය පන්තිය බිහිවීමත්, වැඩිවසම් ක්‍රමය ආරම්භ වීමත්ය. අවුරුදු ගණනාවක් මීනමොනො හා තයිරා යන ක්ෂත්‍රීය පවුල් දෙක අතර පැවති යුද්ධවලින් අවසානයෙහි ජයගත් මීනමොනො යොටිනොමො නම් සේනාපතියා 1185 දී ජපානයේ ප්‍රථම ආඥාදායකයා බවට පත්වී පළමුවැනි ඡෝඟන් ආණ්ඩුව හෙවත් ක්ෂත්‍රීය ඒකාධිපති ආණ්ඩුව (බකුජු) පිහිටුවූයේය.

එයිසයි නම් යනිවරයා ධ්‍යාන බෞද්ධ ධර්ම චීනයෙහි හසළ බුද්ධිමතෙකු බවට පත් වී චීනයෙහි සිට පෙරළා ජපානයට පැමිණියේ මේ කාලයෙහිය. ධ්‍යාන බුද්ධාගමේ දැඩි චිනය හා ආත්ම දමන මාර්ගය සමුරයිවරුන්ගේ සිත් අලවා ගැන්මෙහි සමත් විය. ධ්‍යාන බුද්ධාගම විසින් දර්ශනික මතවාද වලට වඩා ප්‍රායෝගික පක්ෂය වැදගත් කොට සලකන්නට වීමත් එහි උගන්වන මූල ධර්මය ඉතා සරල වීමත් වින්තකයන් හෝ බහුග්‍රාහයන් නොවන සමුරයි වරුන් එය වැළඳ ගැනීමට අතිරේක හේතු වී යයි දැක්විය හැකිය. ඒ කෙසේ වෙතත් එවක පටන් ශතවර්ෂ කිහිපයක් ගත වන තුරු ධ්‍යාන බුද්ධාගමට බලවත්ම රුකුල ලැබුණේ සමුරයි වරුන්ගෙන් බව නොරහසකි. අනතුරුව උදවු මුරොමචි යුගයෙහි (1333-1568) පහළ වූ රාජසභා සංස්කෘතිය සම්පූර්ණයෙන්ම විකාශනය වූයේ ධ්‍යාන බුද්ධාගමේ ආනුභාවයෙන් යයි කියන්නට පිළිවන.

ධ්‍යාන බුද්ධාගම වූ කලි තර්ක ශාස්ත්‍රයට හෝ දුරවබෝධ දර්ශනික මත වාද වලට හෝ ප්‍රිය නොවූ විනුන් හා ජපනුන් විසින් බුද්ධධර්මයේ හරය වටහා ගැන්මට දරන ලද උත්සාහයක ප්‍රතිඵලය වසයෙන් හැඳින්විය හැක.

Figure (b) Letterpress

තුළ පවා ස්වභාව ධර්මය පිළිබඳව බල පැවැති සංකල්පය වූයේ එය පටු කව මාර්ගයක වලයනය වන, සදකාලිකවම විපරිණාම ලෙස පවතින, එසේම, නිව්ටන් උගැන්වූ පරිදි, එහි සදකාලික බගෝල වස්තූද, ලින්තෝ උගැන්වූ පරිදි වෙනස් කළ නොහැකි ජීවේන්ද්‍රීය වර්ගයන්ද අන්තර්ගත කරගත් එකක් බවයි. නවීන ද්‍රව්‍යවාදය, ස්වාභාවික විද්‍යාවේ වඩාත් මෑත සොයාගැනීම් ආලිංගනය කරගනියි. නවීන විද්‍යාවේ සොයාගැනීම් අනුව ස්වභාව ධර්මයාටද කාලයානුබද්ධ ඉතිහාසයක් ඇත. යහපත් තත්වයන් යටතේ බගෝලීය වස්තුවල වෙසෙන ජීවේන්ද්‍රීය දේහයන් මෙන්ම එම වස්තූද උපදිමින්, විනාශවෙමින් පවතින්. ස්වභාව ධර්මය සමස්තයක් වශයෙන් පුනරාවර්තන වක්‍ර මාර්ගයක ගමන් ගනිනැයි කීමට කිසියම් ලෙසකින් තවමත් උවමනා කරනත් එම වක්‍රයන් විශාලත්වයෙන් අනන්තමය පරිමාණයක් ගන්නේය. මේ දෙඅංශයෙන්ම නවීන ද්‍රව්‍යවාදය ප්‍රධාන වශයෙන්ම දයලෙක්තික වේ. එනිසා, රැජිණක් ලෙසට අනෙකුත් විද්‍යාවන් නැමති කළහකාරී රංචුව පාලනය කරන බවත් පෙන්වූ ආකාරයේ දර්ශනවාදයක ආධාරය දැන් එයට උවමනා නැත. මහා ද්‍රව්‍ය සමස්තයක්, ද්‍රව්‍ය පිළිබඳ අපේ දැනුම සම්භාරයක් තුළ විශේෂ විද්‍යාවකට හිමිවන තැන නොවැලැක්විය හැකි ලෙසම පැහැදිලි කිරීමට එයට සිදු වූ වහාම සමස්තය සම්බන්ධයෙන් ව්‍යාවෘතවන විද්‍යාවක් යනු අතිරික්තයක් හෝ අනවශ්‍ය දෙයක් හෝ වන්නේය. මුළු මහත් අතීත දර්ශනවාද සමස්තයෙන්ම තවමත් රැකී පවතින්නේ වින්තන විද්‍යාව සහ එහි නීති, එනම් විධිමත් තර්ක ශාස්ත්‍රය හා දයලෙක්තික ක්‍රමය පමණකි.

කෙසේ හෝ වෙවා, ස්වභාව ධර්මය පිළිබඳ සංකල්පයේ විප්ලවය ඇති කළ හැක්කේ, පර්යේෂණය විසින් සම්පාදනය කරන නියත කරුණුවලට අනුපාත ලෙස වන නමුත්, දැනටමත්, ඉතිහාසය පිළිබඳ සංකල්පය කීරණාත්මක ලෙස වෙනස් වීමට තුඩුදුන් යම් යම් ඓතිහාසික සාධක බොහෝ ඉහතදී ඇතිවී තිබේ. ප්‍රථම කම්කරු පංති කැරැල්ල 1831 දී ලියෝන්හිදී ඇති විය. ප්‍රථම ජාතික කම්කරු පංති ව්‍යාපාරය, එනම් ඉංග්‍රීසි වාට්ස්ව් වරුන්ගේ ව්‍යාපාරය, 1838 හා 1842 අතර කාලයේදී එහි උච්ඡස්ථානයට නැංගේය. යුරෝපයේ දියුණු රටවල් බොහොමයකම ඉතිහාසය තුළ නිර්ධන පංතිය හා ධනපති පංතිය අතර පංති අරගලය පෙරමුණට වන්නේය. එය පෙරමුණට ආයේ එක් අතකින් නවීන කර්මාන්තයේද, අනිත් අනිත් ධනපති පංතිය අළුතින් ලබාගත් දේශපාලන අධිපතිත්වයේද, සංවර්ධනයට අනුපාත ලෙසය. ප්‍රාග්ධනයේ සහ ශ්‍රමයේ අවශ්‍යතා-

Figure (c) Letterpress (Russia)

බුලත්කොහුවිටිය මහරංගල්ල ප්‍රදේශයේ නිවසක අඹුසැමි ආර-
 චුලක් දුරදිග යාමෙන් කෝපයට
 පත් සැමියා දරු දෙදෙනා ද සමග
 පෙරේද (10 ද) නිවෙස අසල
 ගැඹුරු ලිඳකට පැන තිදෙනාම
 මියයාමේ ශෝක ජනක සිද්ධියක්
 බුලත්කොහුවිටිය පොලීසියෙන්
 වාර්තාවේ.

සිද්ධියෙන් පසු බිරිය ද ලිඳට
 පැන ඇතත් අසල්වාසීන් ඇය
 බේරාගෙන ඇතැයි බුලත්කො-
 හුවිටිය පොලීසි ස්ථානාධිපති
 නිමල් ගුණවර්ධන මහතා පැව-
 සිය.

මියගොස් ඇත්තේ බුලත්කො-
 හුවිටිය - මහරංගල්ල පදිංචි
 ආර්. එම්. තරුණ සදිස් රණසිංහ
 (34) (පියා), ආර්. එම්. ඉඳුනිල්
 දිල්ඝාත් (2 1/2) පුතා සහ ආර්.
 එම්. තරුණිකා ප්‍රියංඡිකා (මාස
 10) දියණියයි. ඩබ්ලිව්. ඒ වන්දු-
 ලතා (32) (මව) කැගල්ල මූලික
 රෝහලේ ප්‍රතිකාර ලබයි.

මියගොස් සිටින තරුණ සදිස්
 රණසිංහ මහතා මහරගම රාජික

අධ්‍යාපන ආයතනයේ ලිපික-
 රුවකු ලෙස සේවය කළ අයෙකු
 බවත්, රෝහල් ගතකොට සිටින
 වන්දුලතා මහත්මිය බුලත්කො-
 හුවිටිය ප්‍රදේශයේ පාසලක ගුරු-
 වරියකි.

අඹුසැමි යුවල අතර දින කිහි-
 පයක සිට බහිත් බස් වීම් සිදුවී
 ඇති බවත්, සිද්ධිය වූ දිනයේද
 බිරිය සහ සැමියා අතර ඇති වූ
 බහිත් බස්වීමක් දුරදිග ගොස් ඉන්
 කෝපයට පත් සැමියා දරු
 දෙදෙනා ද සමග අඩි 30 ක්
 පමණ ගැඹුරු ලිඳට පැන තිබේ.

දරු දෙදෙනා ලිඳ තුළදීම මිය
 ගොස් ඇති අතර සැමියා මිය-
 ගොස් ඇත්තේ උදුගොඩ රෝහ-
 ලට ඇතුළත් කිරීමෙන් පසුව බවද
 පැවසේ.

මෙම ශෝකජනක සිද්ධියෙන්
 ප්‍රදේශයේ ජනතාව මහත් කම්-
 පාවට පත්ව සිටින බවත් රෝහල්
 ගතකොට සිටින කාන්තාවගේ
 තත්ත්වය බරපතල නොවන
 බවත් කියයි.

මියගිය අයගේ මහේස්ත්‍රාත්
 පරීක්ෂණය රුවන්වැල්ල වැඩ
 බලන මහේස්ත්‍රාත් සුරච්චර වික්‍ර-
 මනායක මහතා විසින් පවත්වන
 ලදී. එහිදී සාක්ෂි දුන් රංගල්ල

මිට්ට පදිංචි ඩී. පී. සීලවතී මහත්-
 මිය මෙසේ කීවාය. ලිඳට පැන
 කියන ශබ්දයට මම දුටුගෙන
 එනවිට මහත්තයා දරු දෙදෙනත්
 සමග ලිඳට පැනල. කොහේදෝ
 සිට දුටුගෙන ආ තෝතත් ඒ
 වෙලාවේම ලිඳට පැනන.

රුවන්වැල්ල වැඩබලන
 මහේස්ත්‍රාත් සුරච්චර වික්‍රමසිංහ
 මහතා ඉදිරියේ පැවති මරණ
 පරීක්ෂණයේදී තවදුරටත්
 සාක්ෂි දුන් සීලවතී මහත්මිය
 මෙසේ ද පැවසීය.

පසුගිය 10 ද, දවල් 1.00 ට
 පමණ මම ගෙදර සිටියදී වන්දිමා
 දුටු පහලගෙදර පැත්තෙන් කැං-
 හනවා ඇහුණි.

බුට් අසියා ලමසි දෙන්නා සමග
 ලිඳට පතින්න යනවා කියලත්
 කීවා. ගමේ අයත් මේ වෙලාවේ
 දුටුගෙන ආවා. කඹදල බැරිව
 ඉතිමහකින් ලිඳට බැහැල මේ අය
 ගොඩගත්තා.

බිරිය පළමුවත්, දෙවනුව
 පියාත් ගොඩට ගත්තා. දෙන්න-
 ටම කතාකරන්න පුවත්කම
 තිබුණි. ලමසි දෙන්නා මඩ යට
 ගොස් ගිටියා.

ගොඩට ගන්නා විටත් ලම-
 සිත්ට සිහිය තිබුණේ නැ. ගමේ

Figure (d) Dinamina

ලාංකේය කාන්තාවන්ට තමන්ගේ නිවාස තුළවත් බියෙන්
 සැකෙන් තොරව නිදහසේ කාලය ගතකරන්නට නොහැකි
 තරම් බරපතළ තත්ත්වයක් උදාවී තිබෙනවා. වල් වැදුණු
 නීතිය, සදාචාරය පිරිහුණු වත්මන් සමාජය වහාම
 ප්‍රතිශෝධනයට ලක් නොකළහොත් අපේ රටේ හෙට දවසේ
 අනාගතය කෙබඳුවේදැයි සිතා ගැනීමට පවා අපහසුවෙනවා.
 ඇත අතීතයේ අපේ රටේ කාන්තාවට අනුරාධපුරයේ සිට
 මාගම්පුරවරය දක්වා කිසිදු බියකින් සැකයකින් තොරව
 නිදහසේ ගමන්කළ හැකි වුණා. ලංකාවේ කාන්තා
 රැකවරණය ගැන දැන් විවිධ කතිකාවන් පැවැත්වුණත් 1977
 න් පසු බිහිවූ ධාර්මික පාලනයෙන් පසුව අපේ
 කාන්තාවන්ට තනිවම තමන්ගේ ගෙදරවත් බියකින්
 සැකයකින් තොරව ජීවත්වෙන්න බැරිතරමට මේ සමාජය
 දුෂ්ට වී කරුණු නිපුණු වී තිබෙනවා.

අපේ රටේ ජනගහනයෙන් 53% කාන්තාවන් ඔවුන් විවිධවූ
 රැකියාවන්හි නිරත වන්නේ පිරිමින් වගේ හරි හරියට. වැවිලි
 කර්මාන්තයේ රැකියාව කරන ලක්ෂ සංඛ්‍යාත වතු කම්කරු
 කාන්තාවන් ලබන්නේ ඉතාම අල්ප වැටුපක්. ඔවුන් ගෙවන
 දුක්ඛිත ජීවිතය කෙබඳු දැයි අපි දන්නවා.

තවත් ලක්ෂ ගණනක් ඇගලුම් කර්මාන්තයේ නිරත
 වෙනවා. ඔවුන් හොඳට අධ්‍යාපනය ලබපු පිරිසක්. ලක්ෂ
 ගණනක් මැද පෙරදිග ගෘහසේවයේ නියැලී සිටිනවා. මේ
 කණ්ඩායම් දෙක තමයි අපේ රටට අවශ්‍යම විදේශ විනිමය
 උපයාදෙන ප්‍රධාන අංශ දෙක. මැද පෙරදිග ගිහින් දුක් විඳින
 කාන්තාවන් හාමිපුතුන්ගේ දහසකුත් එකක් අතවරයන්ට
 ලක්වෙනවා.

Figure (e) Lankadeepa