

Optical Character Recognition

Software Specification

Version 1.3

Team 3

Quan Chau

Hanchel Cheng

Kevin Duong

Jamie Lee

Ryan Morrison

Eric Rodriguez

Andrew Trinh

Table of Contents

Glossary	1
1: Software Architecture Overview.....	2
1.1 Main data types and structures.....	2
1.2 Major software components.....	2
1.2.1 Model:.....	2
1.2.2 Control:	3
1.2.3 View:	3
1.3 Module interfaces.....	3
1.2.1 Model:.....	3
1.2.2 Control:	4
1.2.3 View:	4
1.4 Overall program control flow.....	5
2: Installation	7
2.1 System requirements.....	7
2.2 Setup and configuration	7
2.3 Uninstalling.....	7
3: Documentation of OCR Modules and Interfaces.....	8
3.1 Detailed description of data structures	8
3.2 Detailed description of pre- and post-processing, and OCR functions and parameters.....	10
Preprocessing.....	10
OCR:	12
Post processing:.....	14
3.3 Detailed description of input images and output text	14
4: Testing Plans	15
5: Development Plan and Timeline	17
6: Branching plan.....	18
Back Matter	20
Index	21

Glossary

Digital image processing (DIP) – functions used to edit the image such as sharpening, rotation, resizing, and etc.

Optimal Character Recognition (OCR) – program that matches characters on an input image to produce a text result

Post-processing – editing of text file produced by the OCR

Pre-processing – editing done to the input image to prepare for OCR program

1: Software Architecture Overview

1.1 Main data types and structures

- Image
- ImageList and ImageEntry
- ObjectHandle
- ViewHandle
- CharProbability
- CharProfile

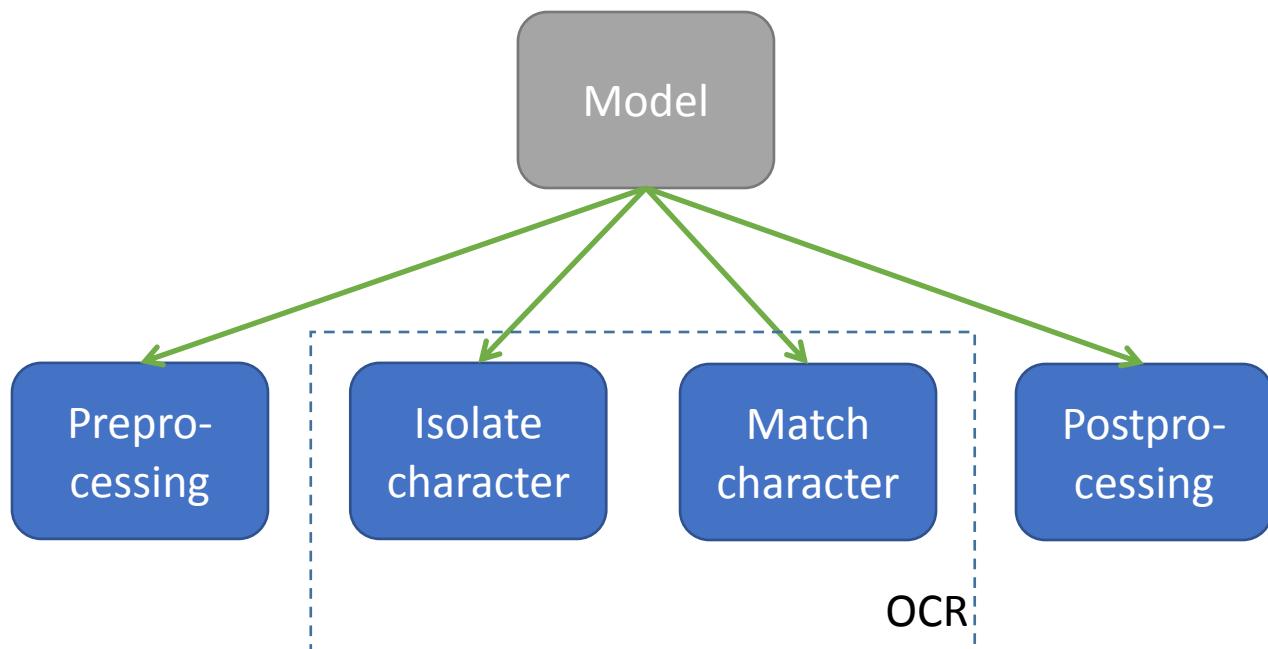
1.2 Major software components

The program will be splitted into three major components: View, Control, and Model.

1.2.1 Model:

Model is responsible for all logic operations and image processing. The model module is divided into four submodules:

- Preprocessing
- Isolate character
- Match character
- Post processing



1.2.2 Control:

Control module processes events, updates the states and call appropriate functions.

1.2.3 View:

View module displays the image and text to the users and receives events initiated from keyboard or mouse by user.

1.3 Module interfaces

1.2.1 Model:

IMAGE * RotatImage (IMAGE * image, int degree);

IMAGE * BlackWhiteImage (IMAGE * image);

IMAGE *StainRemoval(IMAGE *image, int c_var1, int c_var2, int b_threshold, int darken_limiter);

IMAGE *ColorFilter(IMAGE *image, int x, int y, int area_x1, int area_y1, int area_x2, int area_y2, int NewPixelValue, int threshold);

IMAGE *Resize(IMAGE *image, unsigned int newWidth, unsigned int newHeight);

IMAGE *CropImage(IMAGE *image, int x1, int y1, int x2, int y2);

IMAGE *ReadImage(const char *fname);

int SaveImage(const char *fname, IMAGE *image);

IMAGE * PreviewLazyIsolateCharacter(IMAGE * image, FontType Font, int fontsize, int scanres);

ILIST * LazyIsolateCharacter(IMAGE * image, FontType Font, int fontsize, int scanres);

IMAGE * PreviewActiveIsolateCharacter(IMAGE * image, FontType Font, int fontsize, int scanres);

ILIST * ActiveIsolateCharacter(IMAGE * image, FontType Font, int fontsize, int scanres);

UT_string * postProcessing(UT_array *);

UT_array * postProcessingInitializeDictionary();

UT_array * postProcessingInitializeSpecialChar();

UT_string * GetCompileMessage(UT_string * code);

UT_array * IdentifyCharacter(ILIST * imglist, ILIST * Template);

ILIST * InitializeTemplate(void);

1.2.2 Control:

```
ControlHandle * Control_Initialize(void);

void Control_MainLoop(ControlHandle *);

void Control_CleanUp(ControlHandle *);

/*function to take care of all event*/

void Control_ProcessEvent(ObjectHandle * ClickedObject, GdkEvent *);
```

1.2.3 View:

```
GtkWidget *drawMain(ViewHandle * MainViewHandle);

GtkWidget *drawMenuBar(ViewHandle * MainViewHandle);

GtkWidget *drawImageWindow(ViewHandle *MainViewHandle);

GtkWidget *drawTextWindow(ViewHandle *MainViewHandle);

GtkWidget *drawRotateWindow(ViewHandle * MainViewHandle);

GtkWidget *drawCropWindow(ViewHandle * MainViewHandle);

GtkWidget *drawColorFilterWindow(ViewHandle * MainViewHandle);

GtkWidget *drawOCRWindow(ViewHandle * MainViewHandle);

GtkWidget *drawStainRemoveWindow(ViewHandle * MainViewHandle);

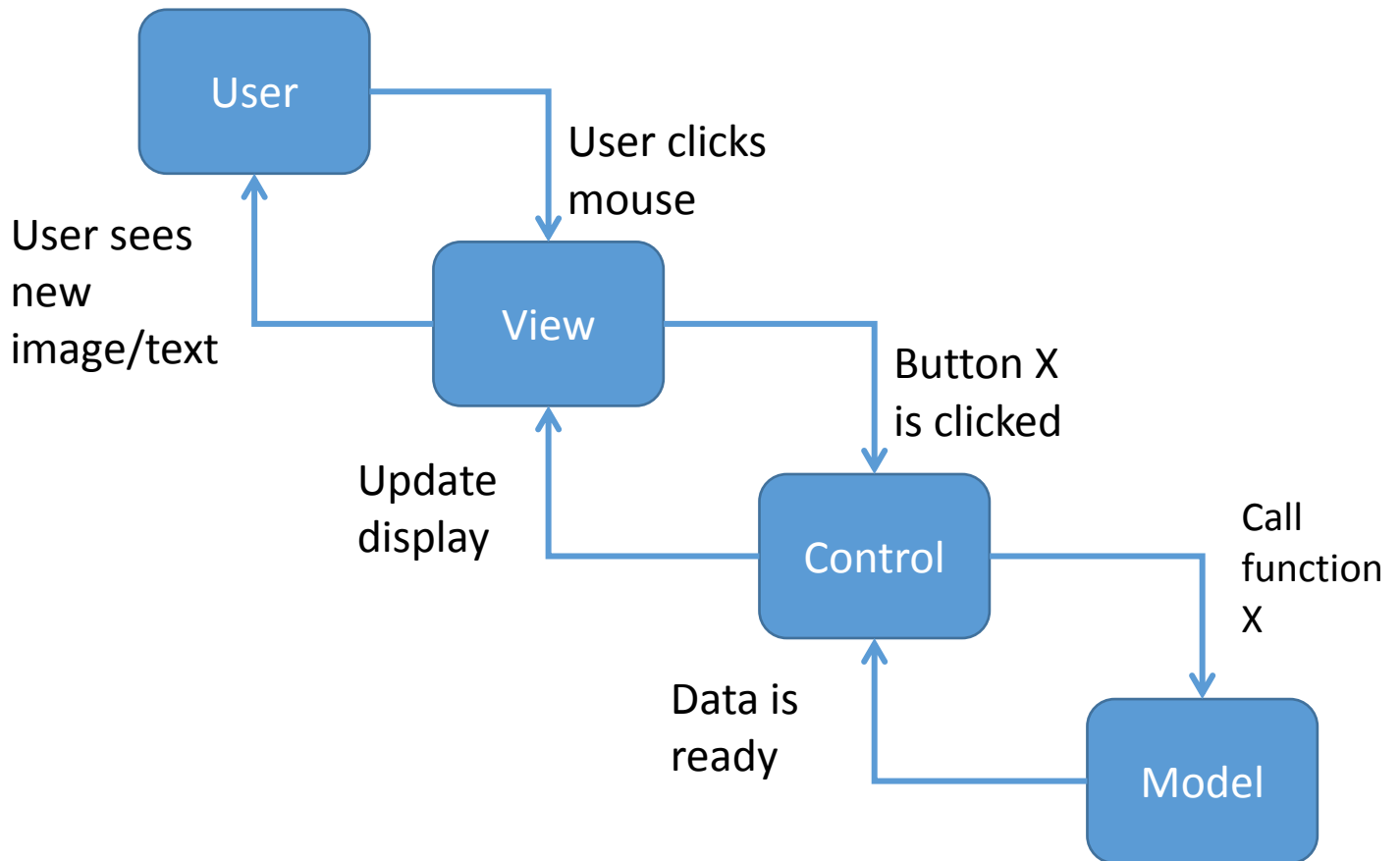
GtkWidget * drawAboutWindow(ViewHandle * MainViewHandle);

GtkWidget *drawHelpWindow(ViewHandle * MainViewHandle);

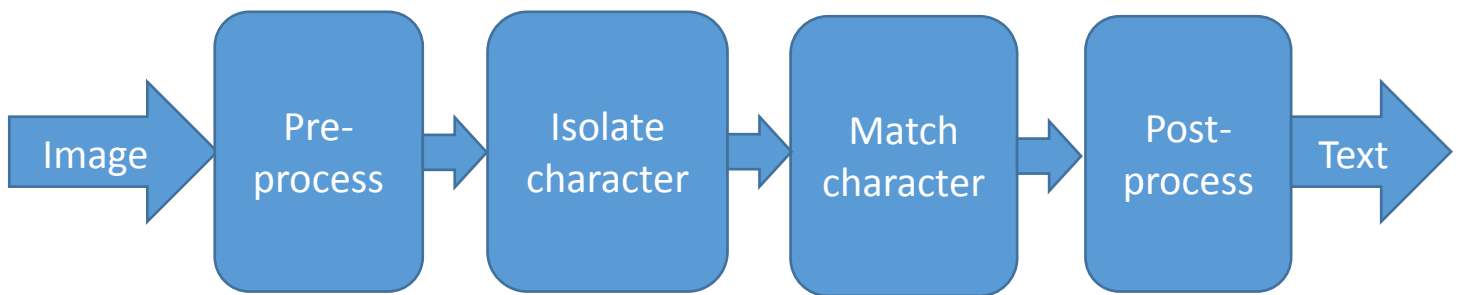
GtkWidget *drawCompileWindow(ViewHandle * MainViewHandle);
```

1.4 Overall program control flow

General view:



Flow chart from an input image to final text result:



2: Installation

2.1 System requirements

- Hardware: PC Hardware (x86_64 server)
- Operating system: Linux OS (RHEL-6-x86_64)
- Dependent third party software:
 - i. gcc
 - ii. GNU make
- Dependent libraries:
 - i. GTK+-2.0 or GTK+-3.0 for graphical user interface
 - ii. Netpbm library for image processing
 - iii. Uthash library for dynamic hash table, array and string

2.2 Setup and configuration

- GTK+ library installation: details are found at www.gtk.org
- Netpbm library installation: details are found at netpbm.sourceforge.net
- Uthash library installation: not necessary because it's already included in the package
- The software comes in a tar.gz package. After downloading, extract the package by running:
`tar -zxvf OCR.tar.gz`
- Change into the directory by running:
`cd OCR`
- Compile the program by running:
`setenv PKG_CONFIG_PATH /usr/share/pkgconfig`
`make OCR`

2.3 Uninstalling

- Change into the directory by running:
`cd OCR`
- Compile the program by running:
`make clean`

3: Documentation of OCR Modules and Interfaces

3.1 Detailed description of data structures

- Image

This data structure hold Red, Green and Blue intensity of a pixel in an image.

```
typedef struct {  
    unsigned int Width, Height;  
    unsigned char *R, *G, *B;  
}IMAGE;
```

This structure has three states:

- Color: R, G and B have distinct values from 0 to 255
- Grayscale: R, G and B have same value from 0 to 255
- Black and white: R, G and B all are 0 (black) or 255 (white)

- ImageList and ImageEntry

These two structures are used to have a dynamic array of Image struct. They are imported from last quarter assignments.

struct ImageList

```
{ int Length;  
  IENTRY *First, *Last;  
};
```

struct ImageEntry

```
{ IENTRY *Next, *Prev;  
  ILIST *List;  
  IMAGE *Image;  
};
```

- ObjectHandle:

This structure holds the name of an object and the GtkWidget * of that object to allow communication between View and Control.

```
typedef struct {  
    char Name[MAX_HASH_KEY_LENGTH];  
    GtkWidget * Widget;  
    UT_hash_handle HashByName;    /*to allow search an object by name*/  
    UT_hash_handle HashByWidget; /*to allow search an object by pointer value*/  
};
```

```
ViewHandle * MainViewHandle;

} ObjectHandle;
```

- ViewHandle:

This structure holds the dictionary (hashtable) of many ObjectHandle. ObjectHandle can be queried by either Name or GtkWidget * value.

```
typedef struct {

    ObjectHandle * ObjectListByName, * ObjectListByWidget;

} ViewHandle;
```

- ControlHandle:

This structure keeps all the data necessary for program to run from functions to functions.

```
typedef struct {

    ViewHandle * MainViewHandle;

    UT_string * InputImageFileName;

    UT_string * OutputFileName;

    ILLIST * MainImageList;

    UT_string * MainOutputString;

    StateEnum State;

    int IsInPreview;

    int InitialImage, InitialText;

    ObjectHandle * XCoordObject, * YCoordObject;

    /*for preview isolate character*/

    IMAGE * PreviewImage;

    /*for resize of image to display nicely*/

    int DisplayWidth, DisplayHeight, OriginalWidth, OriginalHeight;

    /*for post processing*/

    UT_array * Dictionary, * SpecialCharArray;

    /*for identify character*/

    ILLIST * CourierNewTemplate, * LucidaConsoleTemplate;

} ControlHandle;
```

- CharProbability

This structure holds an ascii value and the probability that this ascii value match this single

character image

```
typedef struct {  
    char Char;  
    int Probability;  
} CharProbability;
```

- CharProfile

This structure represents a character spot in the final text string. It contains an array of CharProbability and CharChosen

```
typedef struct {  
    UT_array * CharChoices;  
    char CharChosen;  
} CharProfile;
```

3.2 Detailed description of pre- and post-processing, and OCR functions and parameters

Preprocessing

IMAGE * CropImage(IMAGE *, int x1, int y1, int x2, int y2)

Description: take an image and crop it from coordinate (x1, y1) to (x2, y2). First the function will rearrange two coordinates to get the upper left and lower right corner of the area to crop. Afterwards, it will allocate memory space for new image with new size then set pixel of new image to match those pixels of the cropped area on the original image.

IMAGE * Rotate(IMAGE *, int ClockwiseDegree)

Description: This function will rotate the image clockwise direction by a given degree. The algorithm here is using the pnmrotate from netpbm library with option -background=ff:ff:ff to have a white background.

IMAGE * Resize(IMAGE *, int new_x, int new_y)

Description: This function will resize the image into new size. The algorithm we use is bilinear interpolation. The detailed description can be found here:

http://en.wikipedia.org/wiki/Bilinear_interpolation

IMAGE * ColorFilter(IMAGE *, int x, int y, int area_x1, int area_y1, int area_x2, int area_y2, int NewPixelValue)

Description: This function changes the pixels within a certain rectangular portion of the input image that are similar in color to that of a given reference pixel to some grayscale color (most likely white). The reference pixel may be within any part of the image.

Parameters:

int x & int y:

These parameters define the x and y coordinates of the reference pixel within the image.

int area_x1/x2/y1/y2:

These parameters define the rectangular portion of the image that will be subject to changing.

int NewPixelValue:

This parameter define the grayscale value that pixels (which meet certain conditions) will changed (0 is black, 255 is white, anything between 0-255 is gray).

int threshold:

This parameter represents the level of discrimination when comparing colors to the reference point. Behavior of different threshold values is as follows:

0 ----- 51+

Very little change ----- Entire selection is changed.

This function can be used as an eraser/whitener by setting the threshold to at least 51 and the NewPixelValue to 255.

IMAGE *StainRemoval (IMAGE *image, int c_var1, int c_var2, int b_threshold, int darken_limiter)

Description:

This function is a 3-step process.

The first step is to darken brighter colors in the image. This is to minimize loss of gray or “infected” pixels that may be a part of the actual text.

The second step is to remove (whiten) all pixels whose R, G, B channels differ enough from one another. If the channels do differ enough, they are assumed to be a color which isn't black, gray, or white, and are consequently removed (whitened).

The third step is the same as the second step, except that it only considers pixels that have at least one channel above a certain value.

Parameters:

int c_var1 (Channel Variance 1)

This parameter can be thought of as the coarse adjustment knob. It allows the user to define the minimum value that the R, G, B channels of a pixel must differ from one another in order to be removed (whitened) **when considering the entire image**. int c_var1 can be in the **range of 0 to 51**. The higher the number, the less likely the pixel will be removed (whitened).

int c_var2 (Channel Variance 2)

This parameter can be thought of as the fine adjustment knob. It allows the user to define the minimum value that the R, G, and B channels of a pixel must differ from one another in order to be removed (whitened) **when considering the pixels that have at least one channel above a certain value (i.e. 10 * b_threshold)**. int c_var2 can be in the **range of 0 to 51**. The higher the number, the less likely the pixel will be removed (whitened).

int b_threshold (Brightness Threshold)

This parameter works in conjunction with c_var2. It defines the lower limit of channel (R, G, or B) brightness that a pixel must have in order to be considered for removal. For example, a low b_threshold value will widen the range of RGB values (by including lower RGB values) that are subject to removal. A high b_threshold values will narrow the range of RGB values (by excluding lower RGB values) that are subject to removal. int c_var2 can be in the **range of 0 to 26**. The higher the number, the smaller the range of pixels to be considered for removal.

int darken_limiter

This parameter defines the lower limit of brightness that a single channel (R, G, B) must exceed in order to be darkened. The function uses this value to darken the channels of a pixel which are above (5 * darken_limiter), effectively making brighter colors darker. int darken_limiter can be in the **range of 0 to 51**. The higher the number, the brighter the **channel** of the pixel being considered must be in order to be darkened.

OCR:

ImageList * LazyIsolateCharacter(IMAGE *, FontEnum, FontSizeEnum, ScanResolutionEnum)

Description: This function takes in one IMAGE * and breaks it into multiple IMAGE *, each contains a character. The font type, font size and scan resolution determine how the image is split.

This algorithm works by detecting four margins of the paper by getting the first few black pixels from four sides of the paper, then cut the areas based on pre-defined character width/height and space width/height. This algorithm only works for CourierNew 12pt at 300DPI and LucidaConsole 10pt at 300DPI.

IMAGE * PreviewLazyIsolateCharacter(IMAGE *, FontEnum, int FontSize, int ScanRes)

Description: This function return the original image overlayed with boxes indicating the cropped area of the LazyIsolateCharacters function

ImageList * ActiveIsolateCharacter(IMAGE *, FontEnum, int FontSize, int ScanRes)

Description: This function takes in one IMAGE * and breaks it into multiple IMAGE *, each contains a character or a space that exists between words. The font type, font size and scan resolution determine how the image is split. This algorithm works by detecting the character height based off of the font size and scan resolution. This will determine the maximum number of lines that could be in the cropped text image after a horizontal scan for a black pixel from top to bottom to find the exact y coordinate of the top of the first line. After the exact location of the first line is found, a vertical scan is performed from left to right that goes slice by slice of the length of the character height. A comparison of color is made based off of the character height and if the majority of the pixels in the slice are leaning towards dark shades. This scan produces the start and end boundaries of a character and the start and end of a space between characters and words based off of the aforementioned white/black comparisons. If the whitespace start to end is longer than the width of a character, then the whitespace is pronounced as a space between words instead of characters for the post processing to maintain separation between words. Each character/space is returned as an image in an image list, in addition to a null image for each new line that further allows post processing to recreate the original text.

IMAGE * PreviewActiveIsolateCharacter(IMAGE *, FontEnum, FontSizeEnum, ScanResolutionEnum)

Description: This function returns a copy of the original image overlaid with boxes indicating the cropped area of the ActiveIsolateCharacter function. This is done with similar methodology as the ActiveIsolateCharacter, but instead of cropping and appending to an image list, the box overlay will set the color of pixels in a copy of the cropped image to black, specifically, the pixels that border each isolated character or space.

UT_array * IdentifyCharacter(ILIST *, FontEnum, FontSizeEnum, ScanResolutionEnum)

Description: This function will return an array of CharProfile. Each CharProfile represents a character in the final string and contains all possible characters that this single image matches. The font type, font size and scan resolution determine which font template to select and how the font template is stretched before comparison algorithm works.

The aim of the IdentifyCharacter function is to take in an ILIST from the IsolateCharacter function, and compare each character with a character from a template. The process involves generating a probability for each character in the template to see how it matches up with the ILIST from the IsolateCharacter function. The program's output is a utarray of CharProfiles, which is a structure that contains CharChoices, a utarray full of the CharProbabilities. The IdentifyCharacter function first aligns each character from the template with each character from the ILIST in order to compare the two images. A threshold is used along with an if statement comparing each image pixel

by pixel, and then generating a percentage based on matched pixels. For '\n' and ' ' characters, the IsNewLine and IsSpace functions are used. As for generating the template, an InitializeTemplate function is used.

Post processing:

UT_string * postProcessing(UT_array * CharProfileArray)

Description: This function takes input an array of CharProfile and output a string. The function matches the highest probable character and prints that character into a UT_string and return it.

UT_string * postProcessingAdvance(UT_array *, UT_array *, UT_array *)

Description: This function is the more advanced form of postProcessing(UT_array * CharProfileArray). This function uses an ansi c dictionary and try to match words before a special character (!, @, #, \$, ~, etc) or a space to the ansi c dictionary. If the word matches, it will print that word into a UT_string. Else it will print the highest probable word/characters into the UT_string. If it sees a special character, it will print that into the UT_string itself, without grabbing the next character to form a "word."

3.3 Detailed description of input images and output text

Input images:

This software supports the following input images format:

- PNG
- BMP
- PPM/PGM/PBM
- JPEG/JPG

The conversion from formats other than PPM into PPM uses netpbm library. In particular, the APIs we use from netpbm are:

- pngtoppm
- bmptoppm
- jpegtoppm

After the image is converted into PPM, the file will be read into IMAGE struct defined above and kept in the program memory to be processed

Output text:

Output text will be an array of type char ending with NULL. The char array will contain new line characters as well as normal characters. Here is an example:

H	e	l	l	o	(space)	W	o	r	l	d	\n	E	n	d	\0
---	---	---	---	---	---------	---	---	---	---	---	----	---	---	---	----

4: Testing Plans

Overview:

Module	Function	Input	Output	Description	Phase
<i>Preprocess</i>	Crop	IMAGE	IMAGE	Fully functional	Alpha
	Rotate	IMAGE	IMAGE	Fully functional	Alpha
	Resize	IMAGE	IMAGE	Fully functional	Alpha
	ColorFilter	IMAGE	IMAGE	Fully functional	Alpha
<i>OCR</i>	IsolateCharacter	IMAGE	ILIST	CourierNew, 12pt, 300DPI	Alpha
		IMAGE	ILIST	LucidaConsole, 10pt, 300DPI	Beta
	MatchCharacter	ILIST	UT_array * CharProfileArray	CourierNew, 12pt, 300DPI	Alpha
		ILIST	UT_array * CharProfileArray	LucidaConsole, 10pt, 300DPI	Beta
<i>Postprocessing</i>	PostProcess	UT_array * CharProfileArray	UT_string *	Take only the character with highest probability	Alpha
		UT_array * CharProfileArray	UT_string *	Compare with dictionary	Beta
<i>GUI</i>	Layout + Event	Keyboard + mouse	Event	All operations	Alpha
				Select coordinate by enter number	Alpha
		Keyboard + mouse	Event	Select coordinate by mouse	Beta

Details:

Preprocessing:

- Crop: Take a picture and crop it, must make sure that it takes coordinates from any order
- Rotate: Take a picture and rotate it with any integer degree from 0 to 360. Will test it by issuing random integer degree many times
- Resize: Take a picture and resize it to any new scale. Will test it by issuing random large integer as new scale
- ColorFilter: Take a picture and filter only a certain parts to either black or white using a pre-defined threshold. Will test it by issuing random coordinates

OCR:

- Isolate character: Will test it by reading different images and show the lines to cut into small images. This way is faster since we don't have to save a lot of images and can see the big picture of where the lines are cut

- Identify character: Will take random input from the Isolate Character and try to identify it.

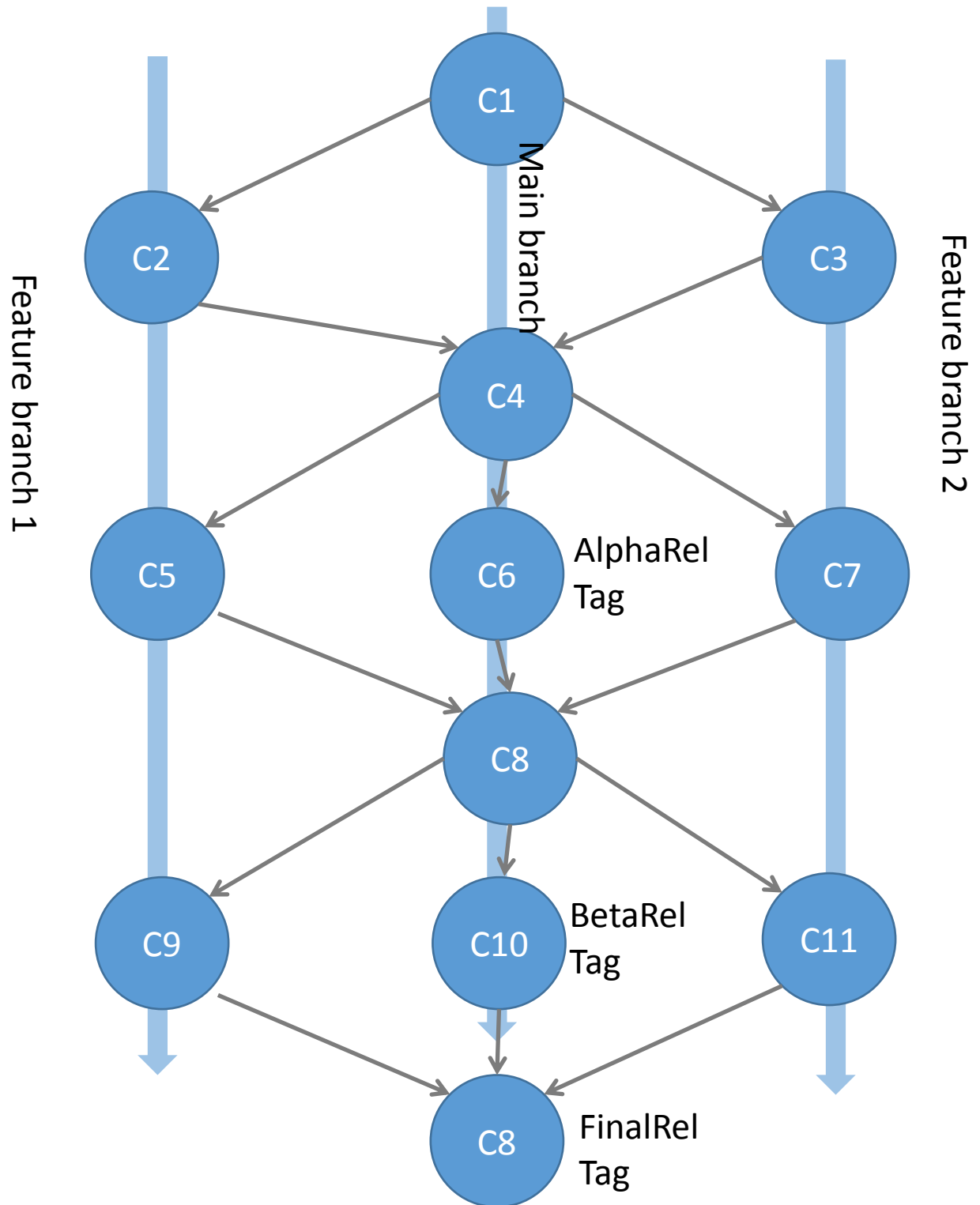
Postprocessing:

- PostProcessing: Generate a random sequence of character probability and takes only character with largest probability. In next phase will try to match the same sequence with as many C keywords as possible.

5: Development Plan and Timeline

	GUI (Ryan + Quan)	Control (Quan)	Preprocessing (Eric)	Isolate character (Hanchel + Jamie)	Identify character (Andrew)	Post processing (Kevin)	Image input (Jamie)
Week1 (by Feb 17)	Pseudo code + try out a real software	Pseudo code + try out a real software	Pseudo code + try out a real software	Pseudo code + try out a real software	Pseudo code + try out a real software	Pseudo code + try out a real software	Pseudo code + try out a real software
Week2 (by Feb 24)	Layout ready	Main flow ready	DIPs + Rotate by any degree	CourierNew, 12pt, 300DPI only	CourierNew, 12pt, 300DPI only	Choose character with highest match	Any file type input
Week3 (by March 3)	Interaction ready	combine many pages, allow external text editor	Solve wrinkle + coffee spot	LucidaConsole, 10pt, 300DPI only	LucidaConsole, 10pt, 300DPI only	Match word with C syntax dictionary	Update if necessary
Week4 (by March 10)	Update if necessary	Update if necessary	Update if necessary	New settings if necessary	New settings if necessary	Run C compiler, make it as compilable as possible	Update if necessary
Week5 (by March 17)	Fix small bug from last release	Fix small bug from last release	Fix small bug from last release	Fix small bug from last release	Fix small bug from last release	Fix small bug from last release	Fix small bug from last release

6: Branching plan



Branching plan explanation:

- 1) C1, C2, C3... represents groups of commits on CVS
- 2) The black arrows coming out from one common commit represents creating a new branch
- 3) The black arrows coming into one common commit represents merging
- 4) Blue arrow represents a feature branch, such as PostProcessing branch and PreProcessing branch

Benefits of this branching plan:

- 1) Avoid pitfall of long-time unmerged branches:
 - The longer the branch stays unmerged, the more chance there are more conflicts to merge
 - This plan forces feature branches to merge frequently (every week) to stay up to date with changes in main branch
- 2) Allow team members to focus on developing features instead of merging
- 3) Merge stable codes only
- 4) Merging and developing are done simultaneously
- 5) Avoid rush to debug new feature to include in release

Back Matter

1. Copyright

This software is licensed under GNU General Public License version 3. Details can be found on <http://www.gnu.org/licenses/gpl.html>

2. Error message:

- a. Image not loaded yet: This error means the user didn't load any image prior to any image operations
- b. Invalid image format: This error means the image selected doesn't have a supported format
- c. Unable to perform OCR: This error means the preprocessed image isn't clear enough for the OCR to detect any characters. To solve this problem, user will have to do some more preprocessing on image before retrying OCR
- d. Invalid input: User entered invalid parameters to images operations

Index

Dependent libraries, 6
Dependent third party software, 6

Hardware, 6
Operating system, 6