# Word Segmentation

-TW015-
Ting-Yu Chang
Advisor: Chi-Cheng Jou
Institute of Electrical Control Engineering
NCTU

July 11th ,2015

# Outline

- **What** is word segmentation?
- **Why** do we need word segmentation?
- **How** do we implement word segmentation?
- Conclusions

# What is word segmentation?

# Introduction-Word

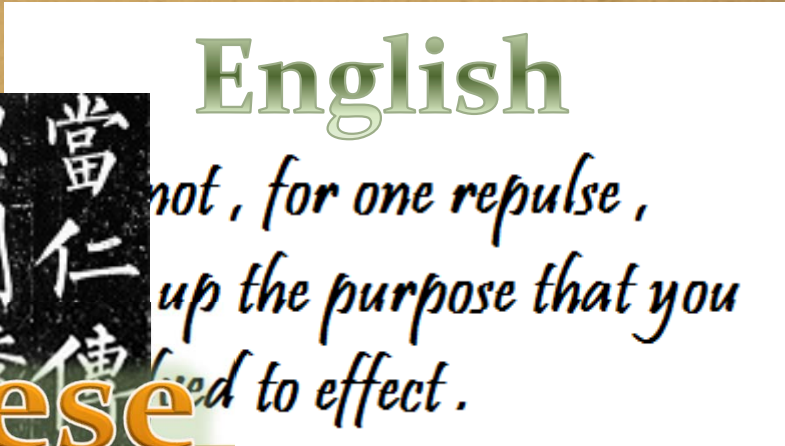**Egyptian**

**Chinese**

**English**
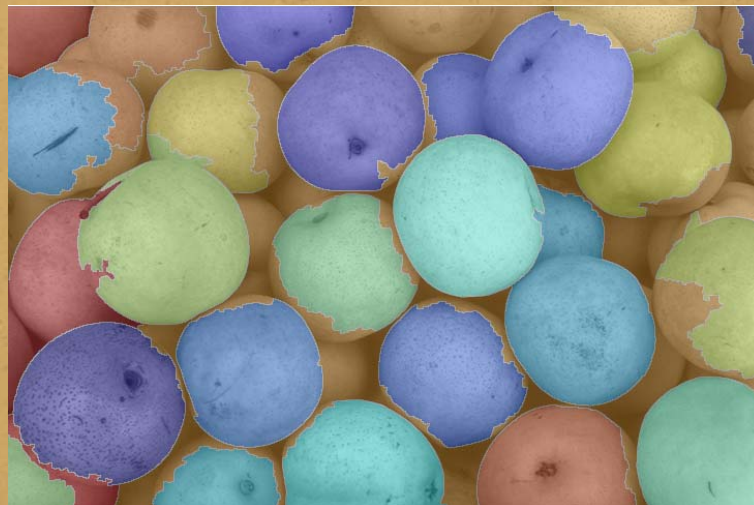
not, for one repulse,

up the purpose that you

...ed to effect.

# Introduction-Segmentation

**Segmentation** is the process of partitioning a digital image into multiple segments



Original Image

Pears
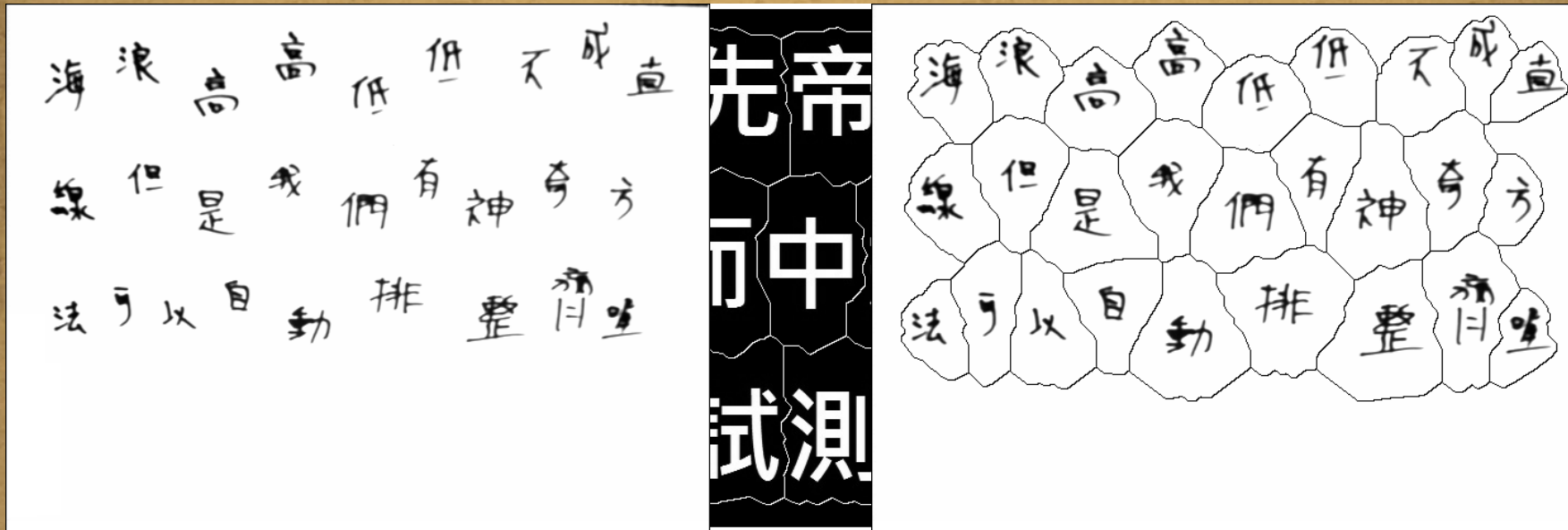
Segmented Image
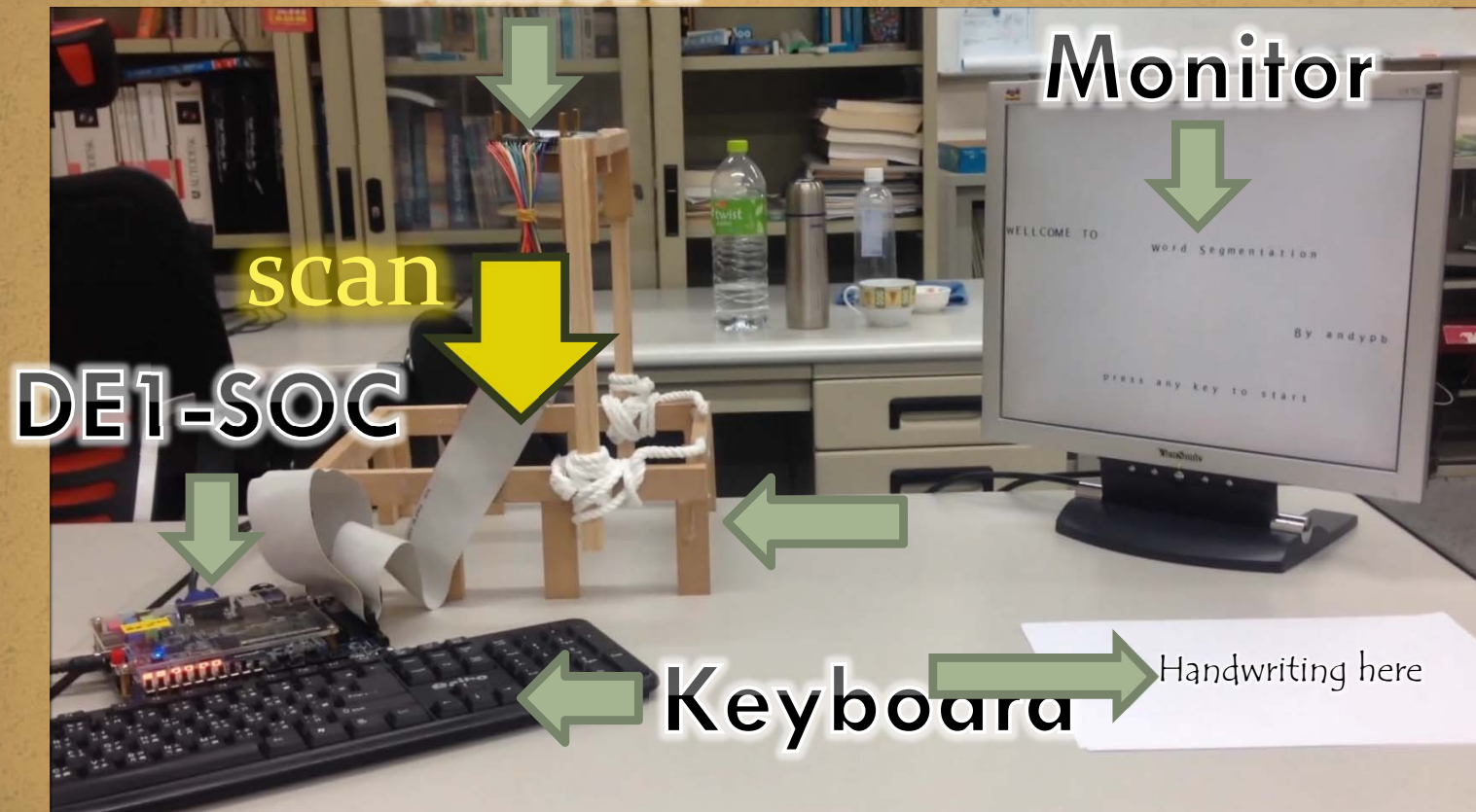
# Word Segmentation



Original Image

Segmented Image

Handwriting words

# Operation Interface

# Why do we need word segmentation?

# Why not......

I can do it!



## Task: counting words

# Yes ! We need it

海浪 高 高 低 低 不 成 直 線 但 是 我 們 有 神 奇 方 法 可 以 自 動 排 整 齊

海浪 高 高 低 低 不 成 直 線 但 是 我 們 有 神 奇 方 法 可 以 自 動 排 整 齊

Delete bounding words

X1000?

How do we implement word segmentation?

# Development Platform

For the software:
For the hardware:

- NIOS II    cpu on FPGA
- HPS(Hard Processor System)    ARM9

TRDB-D5M

DE1-SoC

**Hardware & Software**

Run Histogram equalization
- NIOS II:      seconds
- HPS:      0.065 seconds

**FPGA**

**ARM**

Image processing program

## HPS is Faster!!

# Watershed Segmentation

Watershed segmentation algorithm

Catchments Basin

Catchments Basin

Watershed          Watershed          Watershed

Some words are broken into pieces

# Software Issues

## How to synthesize parts into words?
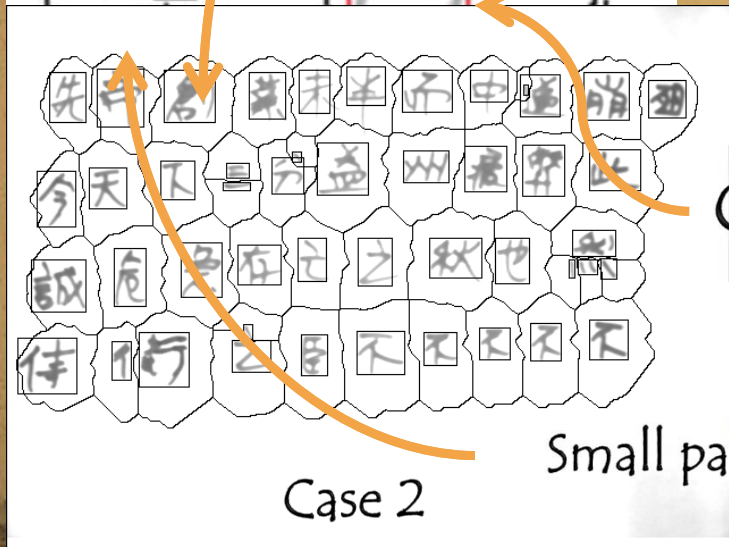
**Solutions:**

a. Identify incorrectly segmented parts

b. Merge the broken parts to each others

A word has 5 parts

# Software Issues(a)



a. Identify incorrectly segmented

Rectangle might be an error

Case 1

Case 2

Small parts might be an error

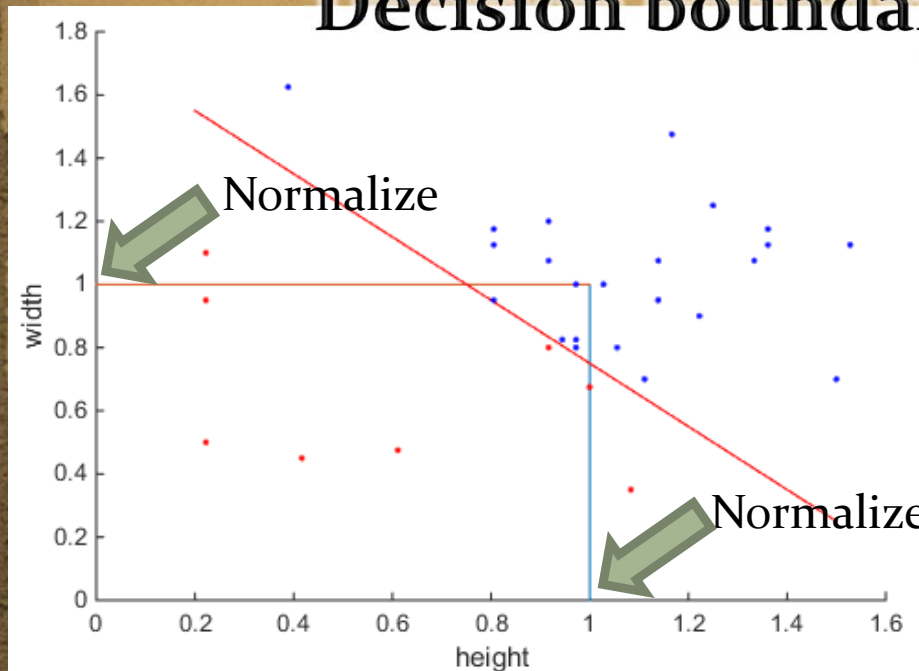**Height-width scatter plots**

Chinese words are square words

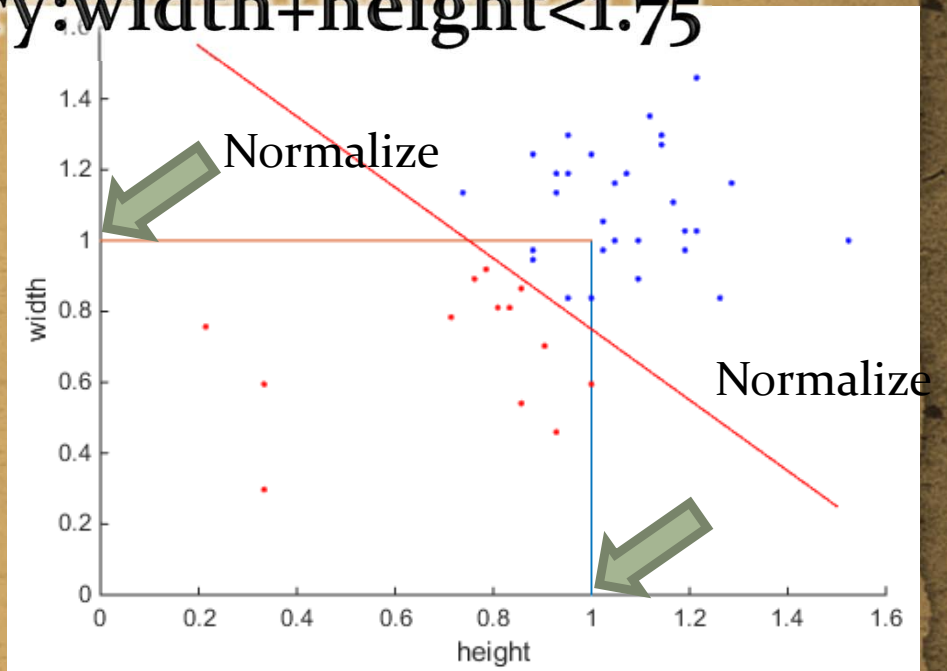# Take a closer look at height-width scatter plot

## Decision boundary:width+height<1.75



Case 1                    Case 2

**Red** point: incorrectly segmented parts
**Blue** point: successfully segmented words
## Robust and Size-invariant

# Software Issues(b)

## b.Merge the broken parts to each others



Merge parts to the nearest neighbor

Calculating minimum/maximum distance between every neighbors

score=0.7minimum distance + 0.3maxmum distance

# Software Issues(b)

## Unfortunetly......

It is too SLOW!!(6 mins per picture)

Mean of the partition

...er picture(...)

Reduce to 10 seconds !!!

6mins → 2mins → 10 seconds

Excellent!!

edge detection

Reduce to only 1 point for each partition

# Conclusions

# Conclusions

**What** is word segmentation?

- For Chinese words
- For hand-writing

**Why** do we need word segmentation?

- For massive documents
- Digital modification: sorting, painting and deleting unwanted word

**How** do we implement word segmentation?

- Use FPGA and ARM
- Hardware issues: capture image data from vga output
- Software issues: merge incorrectly segmented parts
- Performance: 10 seconds per image
- Accuracy: 85% ↑

Thanks for your attention