



mile

MEDICAL INTELLIGENCE AND
LANGUAGE ENGINEERING LAB



MILE lab
Electrical Engineering
Indian Institute of Science (IISc)
Bangalore

MAST

Multi-Script Annotation Toolkit for Scenic Text

T Kasar
Deepak Kumar
M N Anil Prasad
D Girish
Prof. A G Ramakrishnan

- Need for Large annotated data-base
 - Training & bench marking
- Camera-based documents
 - only document image annotation
- Efficient annotation of Scenic text
 - Multi script : flexible adaptation
 - Pixel level ground truth
 - Open source

MAST: Multi-script Annotation toolkit for Scenic Text



WORD LEVEL

LOAD

SEGMENT WORD

THRESHOLD 25



RESEGMENT

RELOAD

BINARIZE

SELECT SCRIPT

NEXT WORD

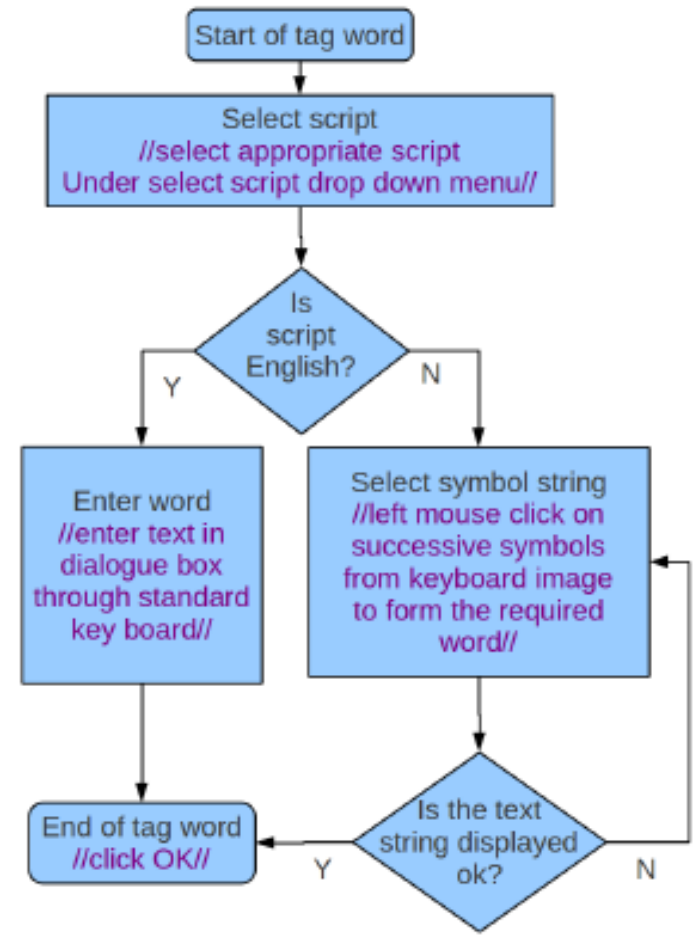
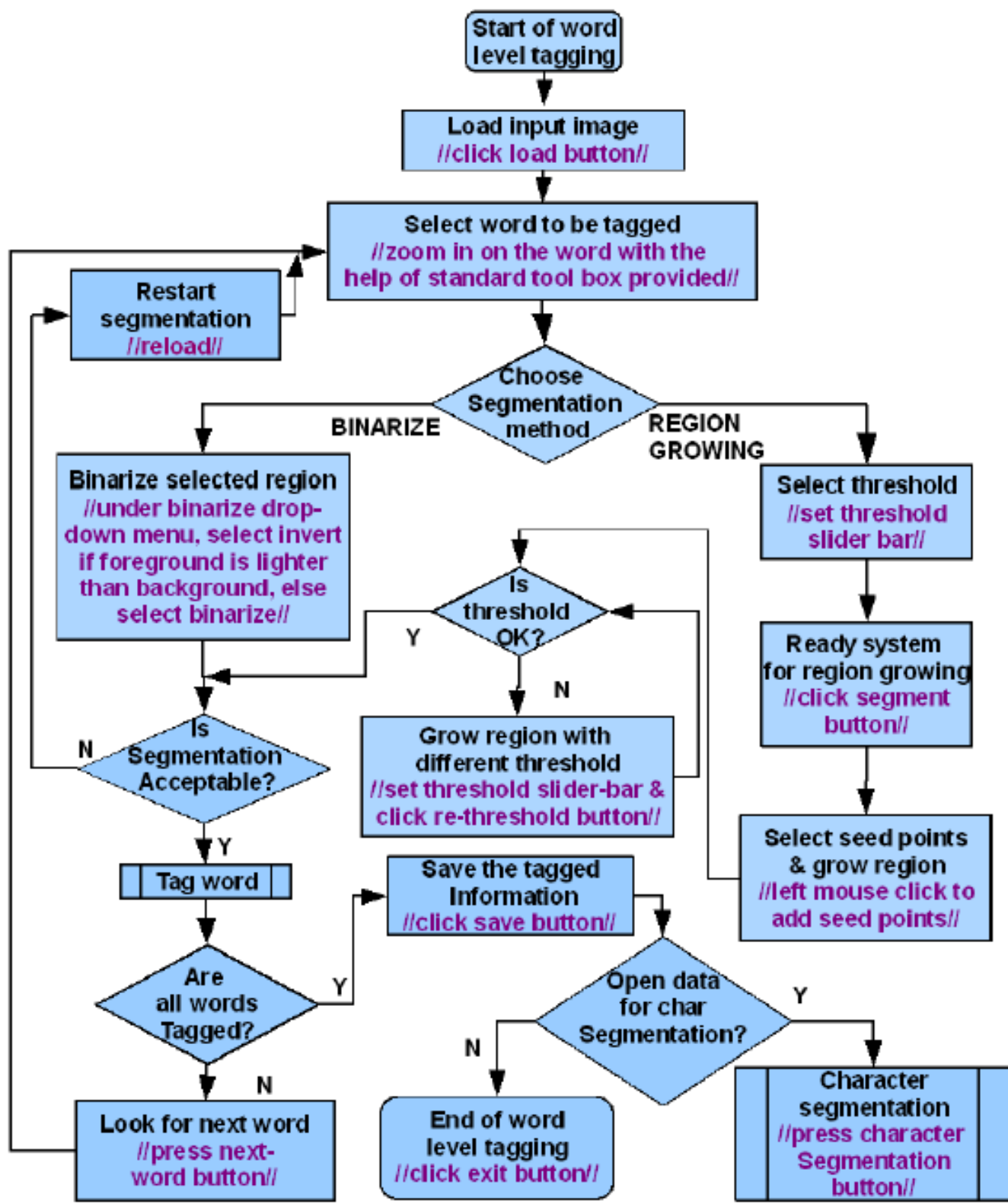
SAVE

WORD BOUNDARY

SEGMENT CHAR

EXIT

- Region growing around manual seed selection
 - CIE L*a*b* color space, with Euclidean distance measure
 - Vary threshold such that pixels corresponding to text are extracted
- If required use Otsu's binarization
- Multiple round :tagging a word at a time



ಆ ರ ೀೋ ಗ ೆಯ

ಆ	ಆ	ಇ	ಈ	ಉ	ಊ	ಋ	ೠ	ಎ	ಏ	ಐ	ಒ	ಓ	ಔ	ಂ	ಃ
್	್	್	್	್	್	್	್	್	್	್	್	್	್	---	---
ಕ	ಖ	ಗ	ಘ	ಙ	ಚ	ಛ	ಜ	ಝ	ಞ	ಟ	ಠ	ಡ	ಢ	ಣ	---
ತ	ಥ	ದ	ಧ	ನ	ಪ	ಫ	ಬ	ಭ	ಮ	---	---	---	---	ಃ	ಽ
ಯ	ರ	ಕ	ಲ	ವ	ಶ	ಷ	ಸ	ಹ	ಳ	ಱ	---	---	಼	ಽ	ಞ
೦	೧	೨	೩	೪	೫	೬	೭	೮	೯	---	---	---	---	✕	∞

OK

UNDO

CANCEL

Ground truth at word level

- `Img.jpg`



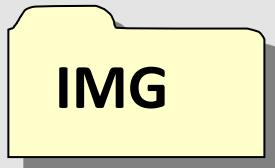
- `Img_seg.jpg`

- `Img.txt`

Bounding box

Script

Word



<code>DSC02617_1.jpg</code>	459 386 195 724	KANNADA	ಆರೋಗ್ಯ
<code>DSC02617_2.jpg</code>	453 1167 169 497	KANNADA	ಕೇಂದ್ರ
<code>DSC02617_3.jpg</code>	826 567 127 560	DEVANAGARI	स्वास्थ्य
<code>DSC02617_4.jpg</code>	762 1195 200 334	DEVANAGARI	केन्द्र
<code>DSC02617_5.jpg</code>	1136 276 127 731	ENGLISH	HEALTH
<code>DSC02617_6.jpg</code>	1127 1080 125 739	ENGLISH	CENTRE

- `Img_1.jpg`



- .
- .
- .

- A Connected component (CC) need not correspond to a character!!!
 - In English 'i'
 - Many Indic languages, a word is a CC



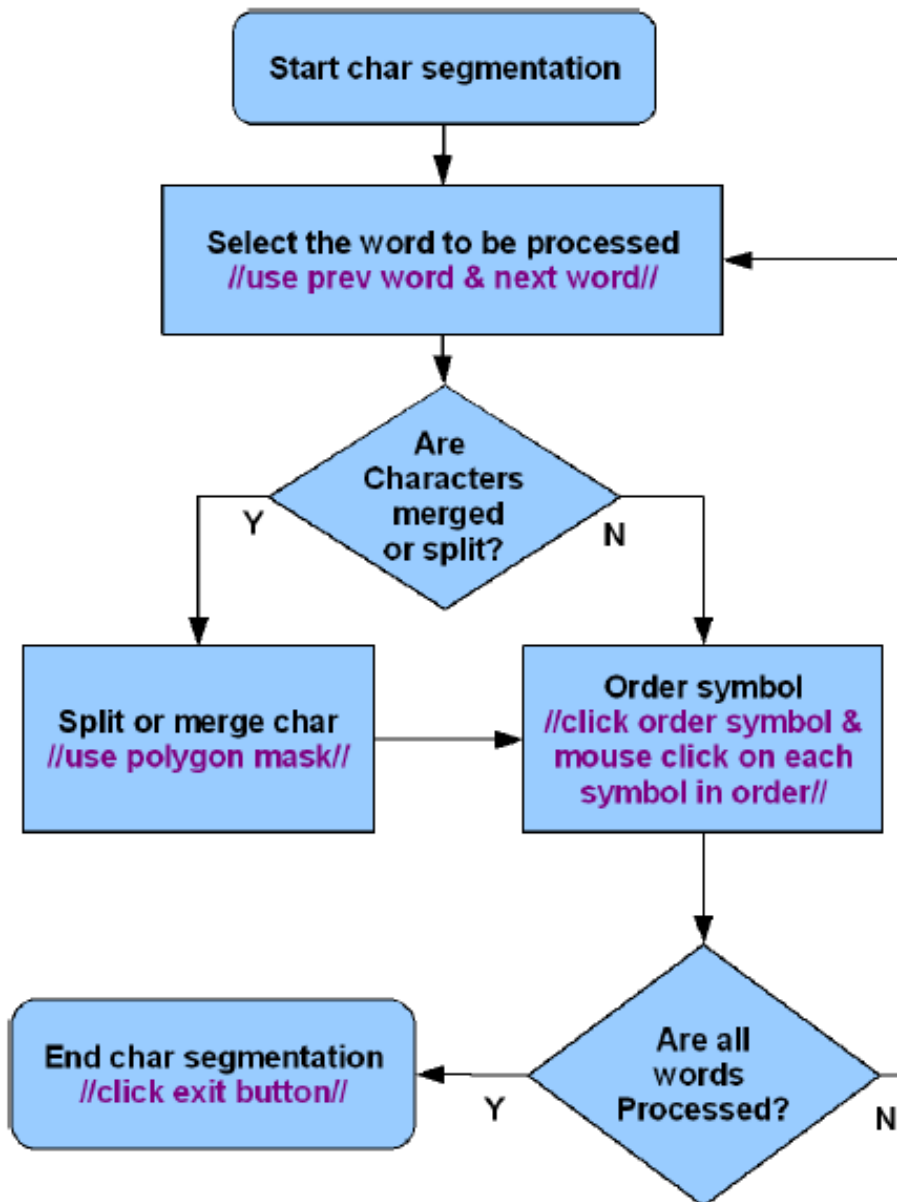
- Therefore need to split or merge CC
- Symbol level as required by the user



mile

MEDICAL INTELLIGENCE AND
LANGUAGE ENGINEERING LAB

Char level





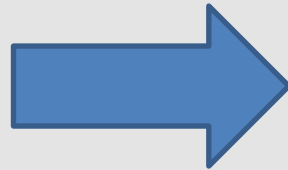
mile

MEDICAL INTELLIGENCE AND
LANGUAGE ENGINEERING LAB

Ground truth at symbol level

IMG

Towna!



IMG

Towna!

Towna!

- Use Keyboard image as i/p to word tagging module
- Perform region growing such that bounding box in the generated text file represents the button
- Tag the buttons Unicode using the standard English keyboard

அ	ஆ	இ	ஈ	உ	ஊ	---	---	எ	ஏ	ஐ	ஓ	ஔ	ஊ	ஃ	஄
஁	ஂ	ஃ	஄	அ	ஆ	---	---	இ	ஈ						
க	ங	ச	ஐ	ஔ	ட	ண	த	ந	ன						
ய	ர	ற	ல	வ	ய	ஷ	ஸ	ஹ	ள						
ஐ	க	உ	ங	ச	ஔ	க	எ	அ	க						


```

tamil_unicode_1.jpg 8 9 51 45 ENGLISH b85
tamil_unicode_2.jpg 8 59 51 50 ENGLISH b86
tamil_unicode_3.jpg 8 114 51 54 ENGLISH b87
tamil_unicode_4.jpg 8 173 51 50 ENGLISH b88
tamil_unicode_5.jpg 8 228 51 50 ENGLISH b89
tamil_unicode_6.jpg 8 283 51 57 ENGLISH b8a
.
.
.
tamil_unicode_61.jpg 234 511 51 56 ENGLISH bf0
tamil_unicode_62.jpg 234 572 51 63 ENGLISH bf1

```

- <http://mile.ee.iisc.ernet.in/mast/>
- MatLab code
- Annotated Database (will be available shortly)
- User guide

- Support for keyboard generation
- Multiple segmentation methodologies
- Semi-automation in tagging text
- Web based