# Deep learning based OCR engine for the Indus script

**Satish Palaniappan**
**(IMSc & Qube Cinemas)**

# Where it all started?

# Computational Epigraphy

Epigraphy is the study of ancient inscriptions, and the place where mathematics and computer science concepts meet epigraphy is called computational epigraphy.

# The Indus Script

- Indus valley / Harappan civilization
- Largest and one of the most ancient civilizations known to mankind
- Major Indus valley sites: (Northwestern regions of India)
  - Harappa, Mohenjo-Daro
  - Chanhu-Daro, Lothal, Kalibangan …
- Inscription's form factor:
  - Stamp seals and Sealings
  - Amulets, stone tablets, pottery …
- The script:
  - Around 3700 text inscriptions according to the M77 corpus, with an average of 5 graphemes per text

# Why Still Undeciphered?

- Paucity of long texts (rarely with 14 graphemes per text)
- Absence of parallel or bilingual text
- No definite knowledge about the underlying language
- The number of graphemes and the very less data (417 symbols with just 3700 texts documented)

# Sample Indus Seals

# Machine Learning in the Study of Indus Script

- In the past, ML has been used for:
  - Classification
    - Based on patterns
    - Based on graphemes
  - Graphemic pattern search
  - Linguistic structure
  - Markov models

# Corpus Formulation

- Need
  - Bottleneck to all ML based research
  - We have been using a 40 years old dataset, with no recent updates
  - Indus script will remain an enigma
- Challenges
  - Extremely laborious human (expert) effort
  - Time consuming to standardize for use
  - Other political issues

# Problem Overview



Output:
89, 17, 8, 342

# Why Deep Learning?

- Complexities in building an OCR engine for reading the Indus script
  - Wear and tear – The artifacts are nearly 4000 year old relics
  - Form factor of the artifacts
  - Very less data with more characters (symbols) to recognize
  - No fixed character set
  - Minute differences between the symbols, account for a completely different representation
  - More undocumented symbols

# Handcrafting the features = Nightmare!



CNN Architecture Diagram

| Stage 1 | Stage 2 | Stage 3 |

Hierarchical Feature Extraction in stages

Image: Pixels → 1st Stage: Edges → 2nd Stage: Object Parts → 3rd Stage: Objects

# Transfer Learning and Fine Tuning

# Data Augmentation

- To compensate for the meager data, we devised certain image augmentation techniques in addition to what Keras has.
  - Vertical and horizontal flips
  - Shear, Crop, Swirl
  - Rotate, Scale, Translate
  - Randomised artificial lighting

# Datasets Used

**Text-NoText Dataset**



Text

No-Text

Both

**Symbols Dataset**



'Jar' sign present

'Jar' sign absent

# Base CNN Architectures Used

**GoogLeNet**



**SymbolNet**

# The Pipeline

# Region Proposal

Proposes regions that have a high probability of containing a symbol, animal, deity, or any iconographic element.



Indus Seal Scan — Extract Seal — Selective Search — Region Grouping

# Extract Seal

Removes the irrelevant background information from the input artifact image.

- Steps
  - Grayscaled and smoothened using multi-scale gaussian filter
  - Threshold the image at the mean pixel value of the background
  - Optimized canny edge detection

# Selective Search

Proposes an array of all possible regions, likely to hold Indus script symbols, depictions of animals, etc.

- Combines the advantages of exhaustive search and segmentation
- Hierarchical grouping of the region proposals based on color, texture, size and fills
- Grid search over the four Selective Search parameters to use the Scale, Min Size, Min Area, and Sigma

# Region Grouping

Improvises the quality of the regions proposed by Selective Search (four-level region grouping and filtering hierarchy)

- Merge concentric proposals
- Contained boxes removal
- Draw super box
- Draw extended super box

The last two levels leverage the prior information that the text regions are contiguous, being mostly arranged along a same line or axis rather than randomly distributed in space.

# Text Region Extraction

Produces exact text-only regions by eliminating the non-graphemic parts off of the region proposals in hand.



**Selective Search**

**Region Grouping**

**Region Classification**

'text'  'both'

**Text Region Formulation**

# Region Classification

All the regions are classified into types: "Text", "No-Text" or "Both".



**ImageNet weights - Frozen (Transfer Learnt)**

**ImageNet weights - Initialized**

Inception Module

**Learning Rate - 2X (Fine-tuned)**

# The Region Classification CNN's Results

**Top-1 Accuracy Scores**

| GoogLeNet's Levels | Top1 Accuracy Scores |
| --- | --- |
| Level 1 (1/3rd network depth) | 87.14% |
| Level 2 (2/3rd network depth) | 87.86% |
| Level 3 (full network depth) | 89.30% |



**Expected Label:** Both
**Actual Label:** Both

**Expected Label:** Text
**Actual Label:** Text

**Expected Label:** No Text
**Actual Label:** No Text

**Expected Label:** Text
**Actual Label:** Text

**Expected Label:** Both
**Actual Label:** Text

**Expected Label:** Both
**Actual Label:** No Text

# The Region Classification CNN's Graphs



(a)   (b)   (c)

# Text Region Formulation

Builds the text-only regions from the labeled region proposals (two level hierarchy)

- Draw TextBox
  - Merge Two "text" regions or a "text" region and a "both" region that are aligned along the same horizontal or vertical axis => TextBox
- Trim TextBox
  - Clip off the non-textual information ("no-text") in those pairs of region proposals, where a "text box"/"text" region and a "no-text" region were overlapping

# Symbol Segmentation

Segments out the individual graphemes from the precise text-only region proposals.



**Text Region**

1.Gray Scale
2.Otsu Thresholding
3.Gaussian Blur

4.Connected Colour Components
5.Combine Component Regions
6.Crop ROIs

# Symbol Segmentation – The Algorithm

---

**Algorithm 1** Symbol Segmentation Algorithm

---

1: **procedure** Segment–Symbol(Image I)
2:     Gray_Image = **Gray_Scale**(I)
3:     Thresholded_I = **Otsu_Thresholding**(Gray_Image)
4:     Smoothened_Image = **Gaussian_Blur**(Thresholded_Image)
5:     Component_ROIs = **Connected_Colour_Components**(Smoothened_Image)
6:     ROIs = **Combine**(Component_ROIs)
7:         Unique_ROIs = **Contained_Boxes_Removal**(Component_ROIs)
8:         Super_ROIs = **Draw_Super_Box**(Unique_ROIs)
9:         ROIs = **Draw_Extended_Super_Box**(Super_ROIs)
10:     Segmented_Symbols = **Crop**(ROIs, I)
11: **end procedure**

# Symbol Identification

Takes individually cropped graphemes from the previous stage and classifies them into one of the 417 symbols (M77 Corpus)



Symbol-wise Segmented Text Region

Jar / No-Jar Classifier

No-Jar    No-Jar    No-Jar    Jar

Classified Symbols

# Symbol Identification - CNN Architecture

- Detects the presence or absence of the "Jar" sign - Binary classifier
- No transfer learning
- New architecture trained from scratch
- Accuracy score of **92.07%** when evaluated over the validation set of the "Jar-NoJar Dataset"

| Data (32*32*3) | → | Convolution 5x5 (Stride 1, 20 Outputs) | → | Convolution 5x5 (Stride 1, 50 Outputs) | → | Dropout | → | Fully Connected (500 outputs) | → | ReLU (Non Linearity) | → | Fully Connected (2 Outputs) | → | SoftMax Classifier |

# Symbol Identification CNN's Graphs



(a)    (b)    (c)

# An Example Flow



Input Image

Region Proposal
- Extract Seal
- Selective Search
- Region Grouping

Region Classification
- 'no-text'
- 'text'

Text Region Extraction
- Text Region Formulation

Symbol Segmentation

Symbol Identification
- Sign 89
- Sign 17
- Sign 8
- Sign 342

# Evaluating the Pipeline

| Stages in Pipeline | Output Classes | | | | | | Indicative Accuracies (completely perfect cases only) |
|---|---|---|---|---|---|---|---|
| Region Proposal and Text Region Extraction | Full Text regions | | | Partial Text regions | | | Full Text Regions (43/50) |
| | 43 | | | 7 | | | 86% |
| Symbol Segmentation | Full Symbols | Partial/ Combined Symbols | No Symbols | Full Symbols | Partial/ Combined Symbols | No Symbols | Full Symbols ((29+5)/50) |
| | 29 | 11 | 3 | 5 | 2 | 0 | 68% |

# Limitations

- Non ML based techniques used in Region Proposal and Symbol Segmentation stages
  - Leads to some parts of the text region and symbols getting missed
- Extremely damaged seals can never be fully read
- Very complex grapheme arrangement structure
- Limited to identify only few frequent symbols of the 417
  - Issue is with the availability of such a dataset, it needs to compiled

# Generalizing to Other Ancient Inscriptions

- We need
  - Sufficient amounts of data in the specified format for the new inscription
  - Tweak the region grouping and symbol segmentation stages and their parameters to harness the new inscription's ergonomics

# The Buzz !

**CIPHER WAR**

*After a century of failing to crack an ancient script, linguists turn to machines*

by Mallory Locklear | Jan 25, 2017, 9:00am EST

Meanwhile, Ronojoy Adhikari, a physics professor at The Institute of Mathematical Sciences in Chennai, India, and his research associate Satish Palaniappan are working on a program that can accurately extract symbols from a photo of an Indus artifact. "If an archaeologist goes to an Indus site and finds a new seal, it takes a lot of time for those seals to actually be mapped and added to a database if it's done manually," says Palaniappan. "In our case the ultimate aim is just with a photograph of a particular seal to be able to extract out the text regions automatically." He and Adhikari are working on building an app that archaeologists can bring to a site on a mobile device that will extract new inscriptions instantly.

SBS

5 MAR 2017 · 8:26PM

## An app to Decipher ancient symbols

PODCAST
tamil_170305_641975.mp3

A professor from Institute of Mathematical Sciences (IMSc), Chennai, and an engineer have developed an app that will allow archaeologists and amateur history buffs alike to, say, capture images of seals on pottery and share it online via the app to assist experts devoted to the recognition and transcription of the script. It will also provide an approximate date by recognition of the iconography and its style... Satish Palaniappan, an SSN College of Engineering graduate who worked on the app talks to Kulasegaram Sanchayan about it.

## Researchers Look To Widen Script Database, Solve Mystery
# New app could help decipher ancient Indus Valley symbols

U.Tejonmayam
@timesgroup.com

**DAWN OF TIME** — The computer application can be used to identify elements belonging to the Indus script

# Chennai team taps AI to read Indus Script

The algorithm uses 'deep neural networks' which are also used in self-driving cars

SHUBASHREE DESIKAN
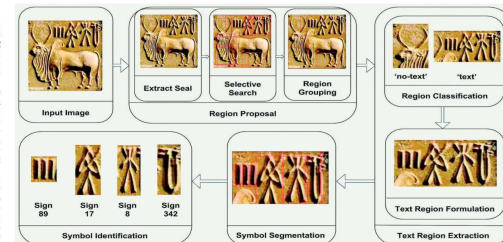
**Step by step:** Scanning the image, the algorithm smartly processes the data in three steps to place its elements within the standard corpus. SPECIAL ARRANGEMENT

**Dr Mahadevan says, "It [the algorithm] represents a significant advance in the computerised study of the Indus Script. I wish I had this software 40 years ago when I compiled the Indus concordance."**

# Thank You!

Satish Palaniappan