

# 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
  - Chanhу-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



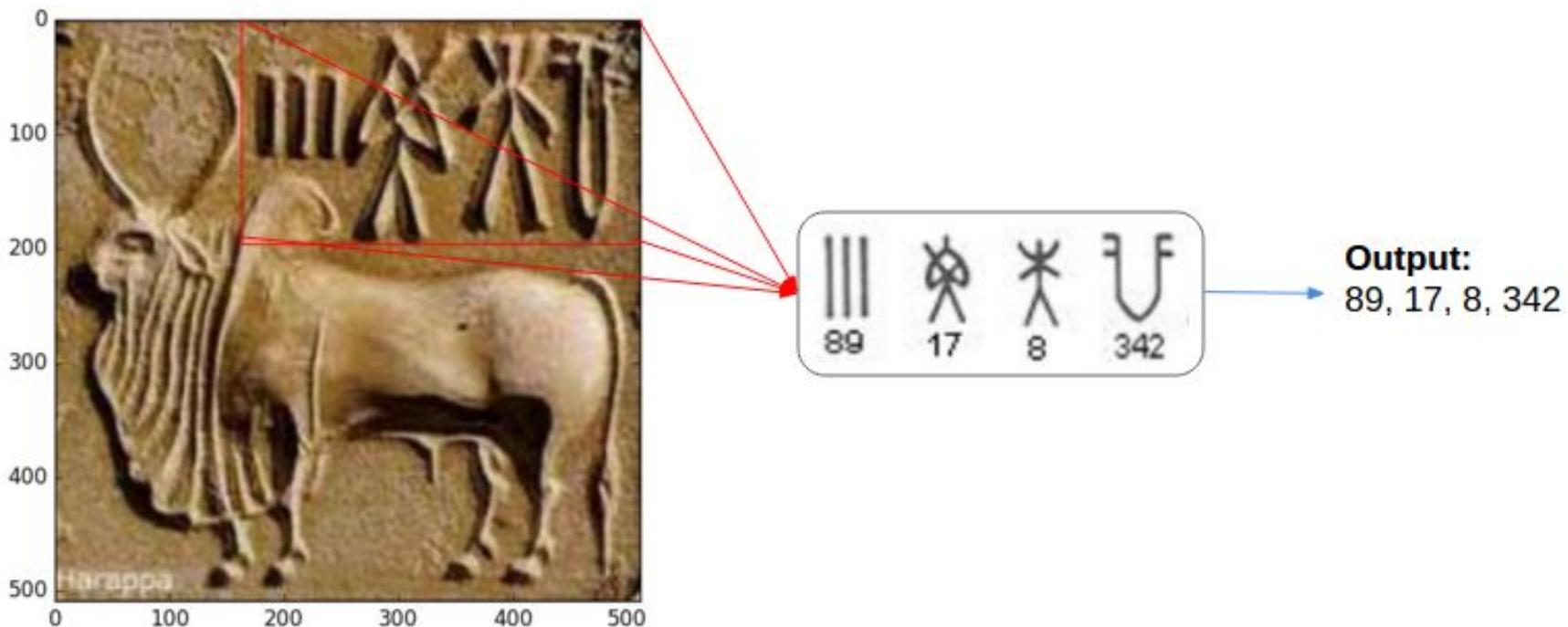
# MACHINE LEARNING IN THE STUDY OF INDUS SCRIPT

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

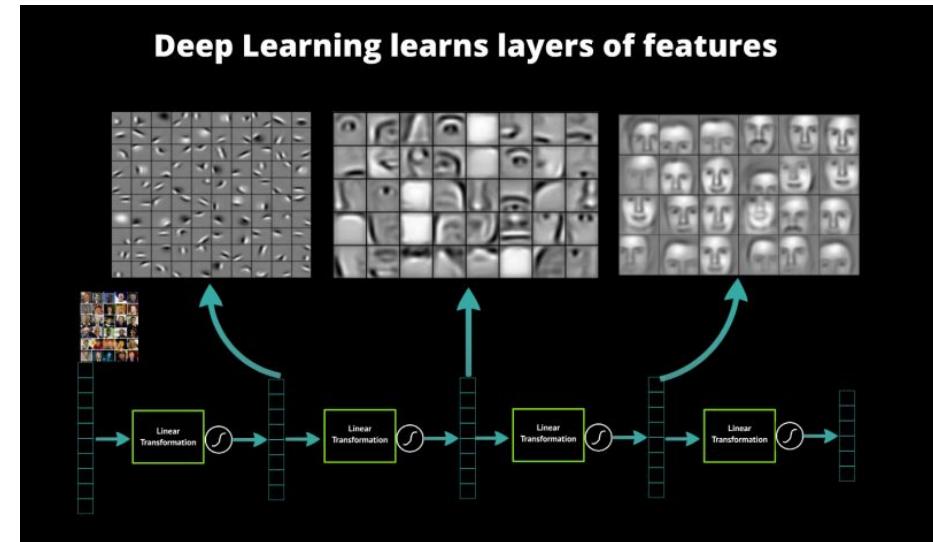
# CORPUS FORMULATION

- 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



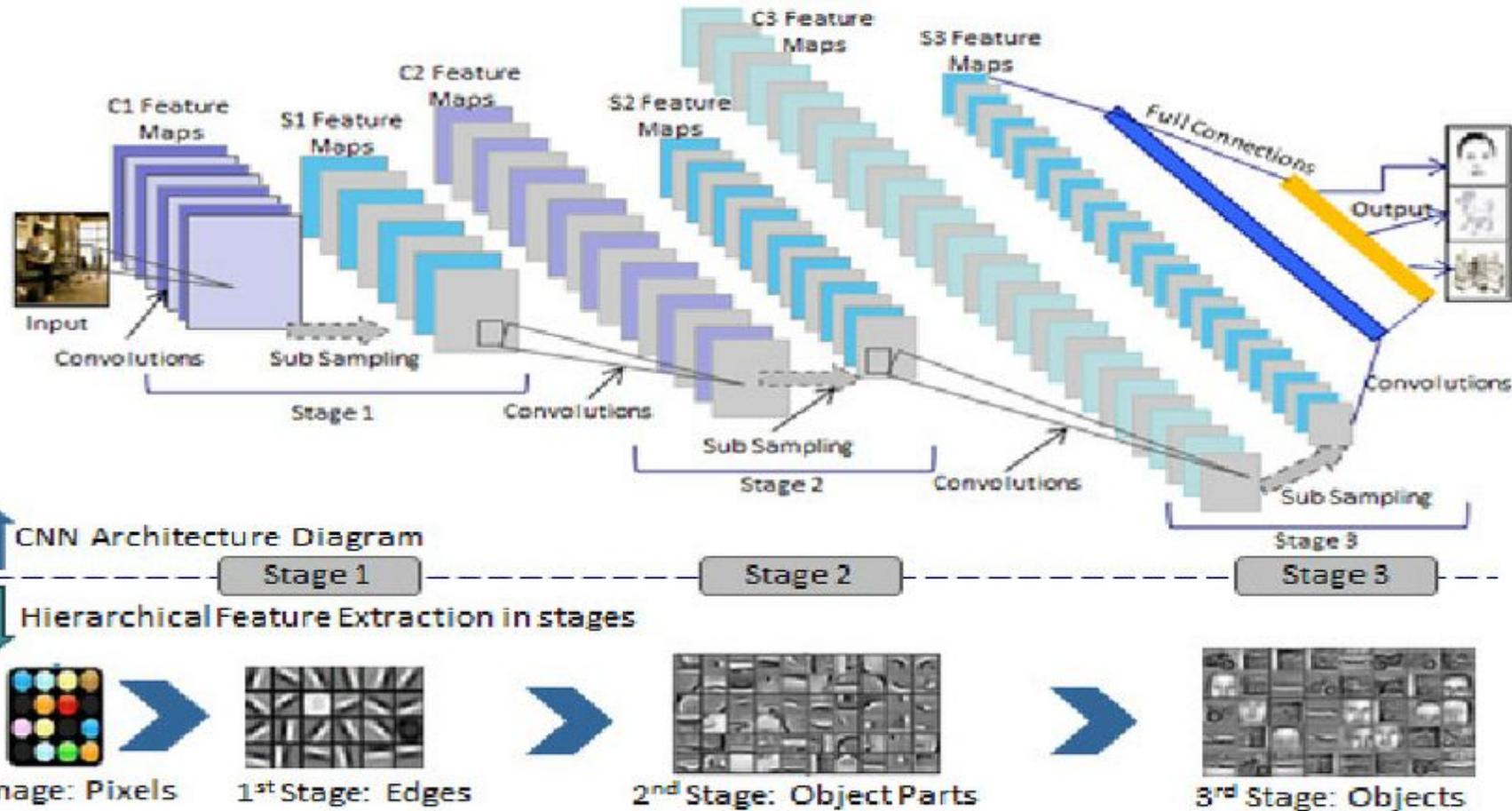
# WHAT IS DEEP LEARNING?



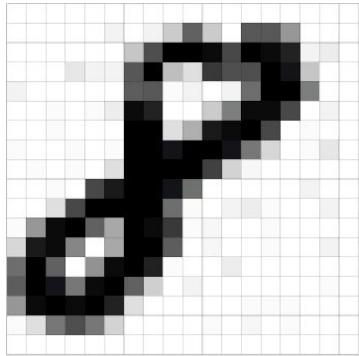
# WHY DEEP LEARNING?

- 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!



# CNNs



1 <small><math>\times 1</math></small>	1 <small><math>\times 0</math></small>	1 <small><math>\times 1</math></small>	0	0
0 <small><math>\times 0</math></small>	1 <small><math>\times 1</math></small>	1 <small><math>\times 0</math></small>	1	0
0 <small><math>\times 1</math></small>	0 <small><math>\times 0</math></small>	1 <small><math>\times 1</math></small>	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

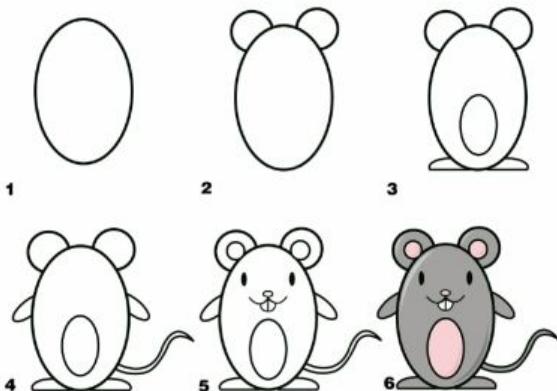
Convolved Feature



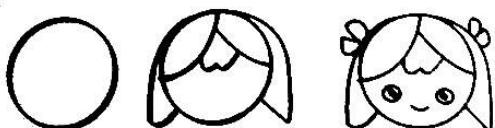
Input

# WHAT IS TRANSFER LEARNING?

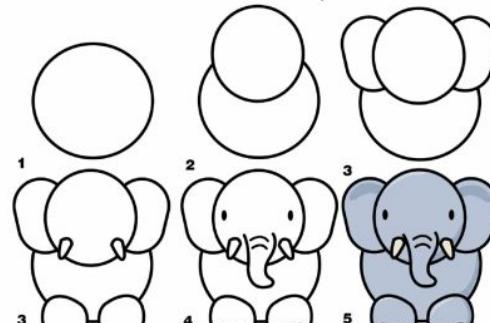
Learn to Draw a Rat



Learn to draw with [www.ActivityVillage.co.uk](http://www.ActivityVillage.co.uk) - Keeping Kids Busy

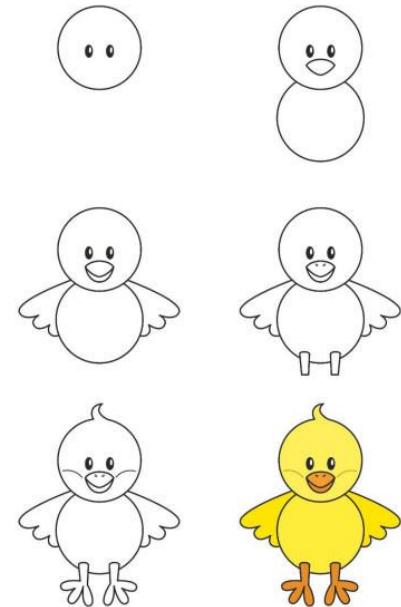


Learn to Draw an Elephant



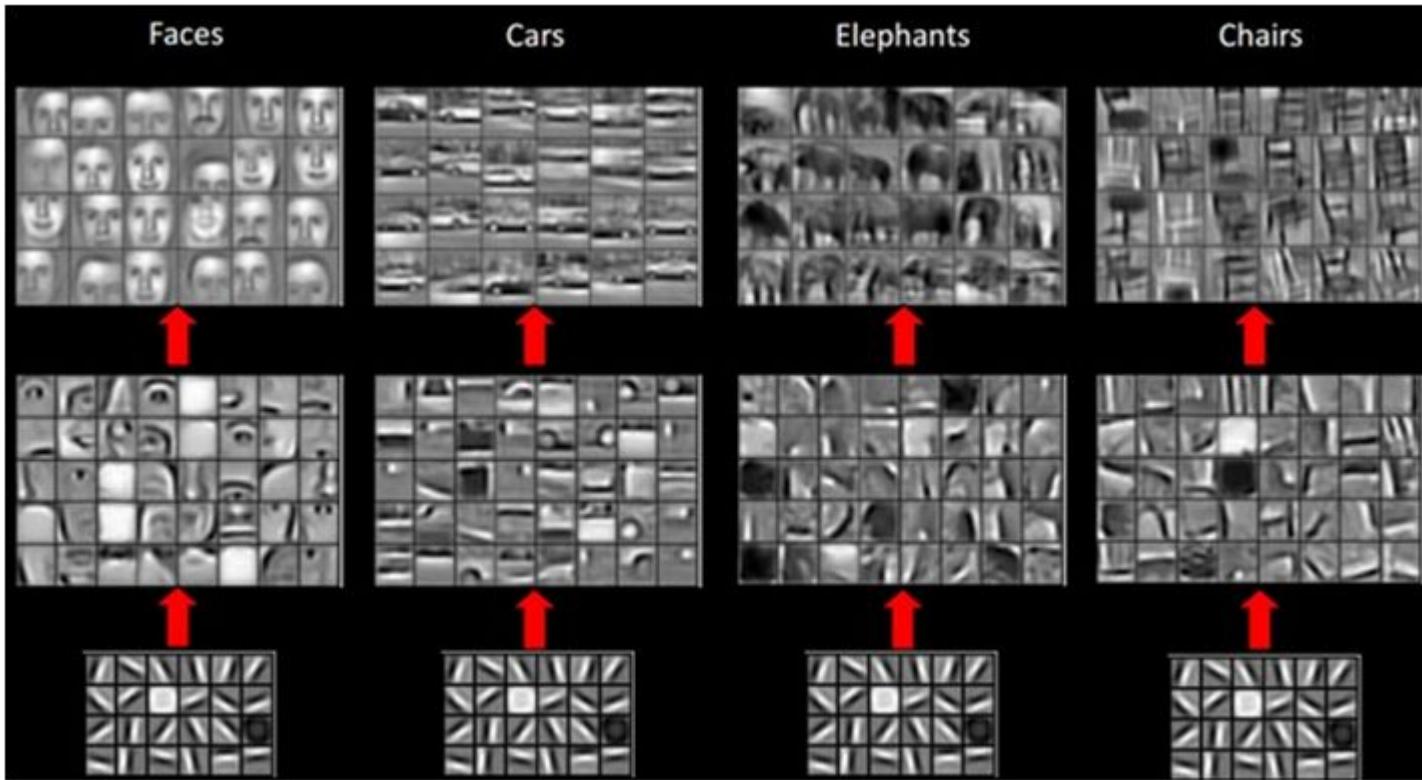
Learn to draw with [www.ActivityVillage.co.uk](http://www.ActivityVillage.co.uk) - Keeping Kids Busy

Learn to Draw a Chick



Learn to draw with [www.ActivityVillage.co.uk](http://www.ActivityVillage.co.uk) - Keeping Kids Busy

# 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

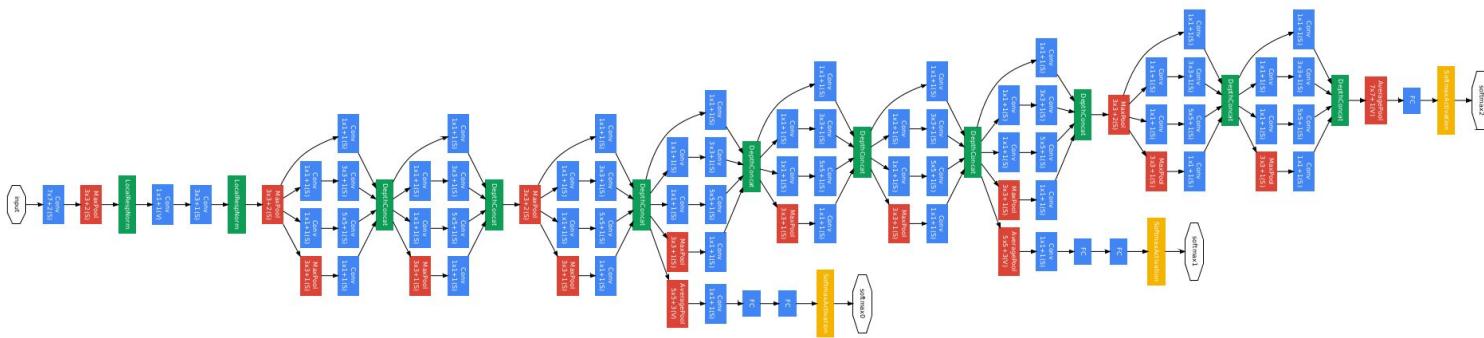


'Jar' sign absent

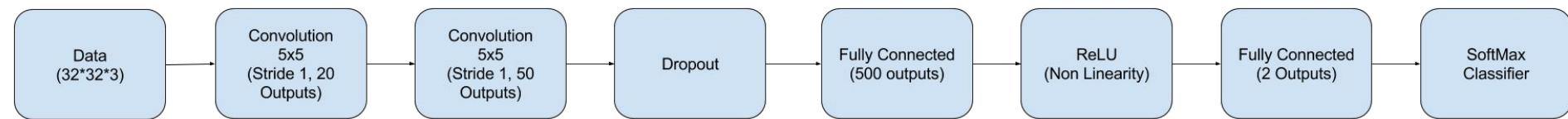


# BASE CNN ARCHITECTURES USED

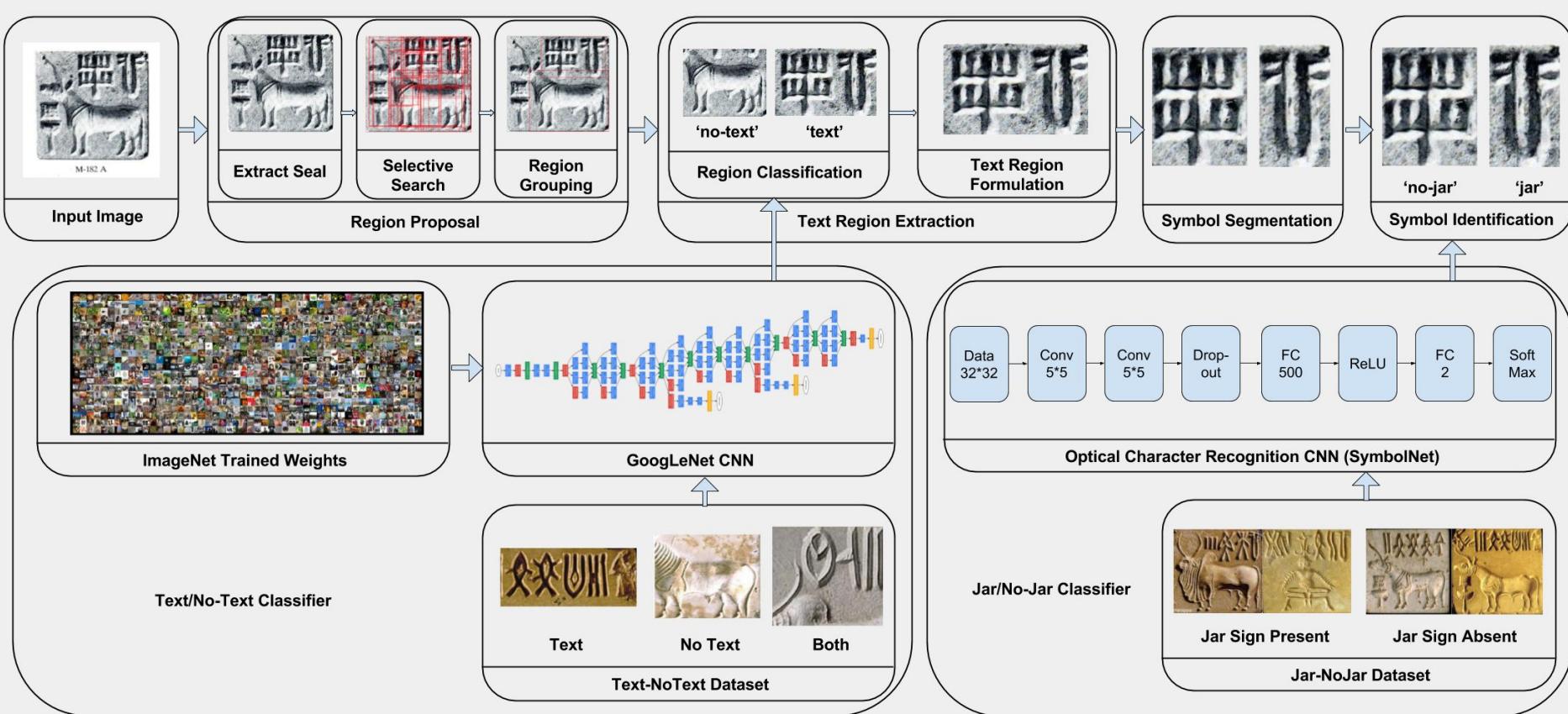
## GoogLeNet



## SymbolNet

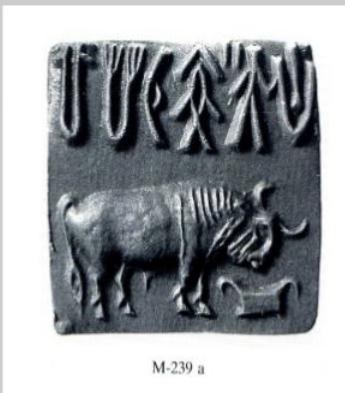


# THE PIPELINE

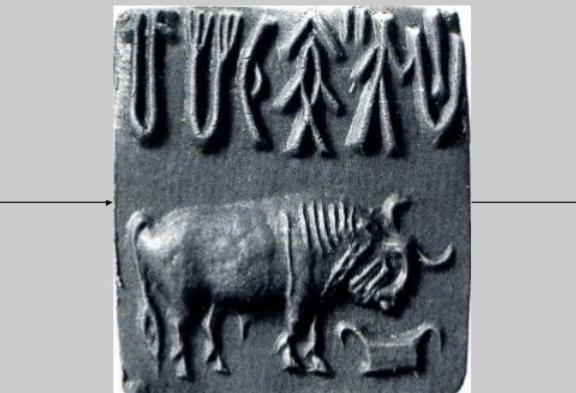


# REGION PROPOSAL

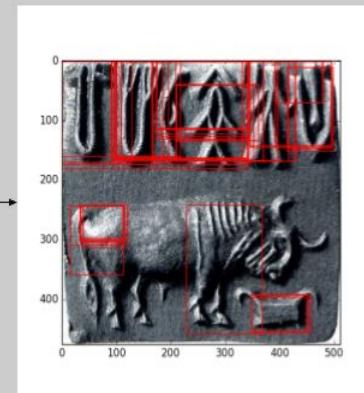
Proposes regions that have a high probability of containing a symbol, animal, deity, or any iconographic element.



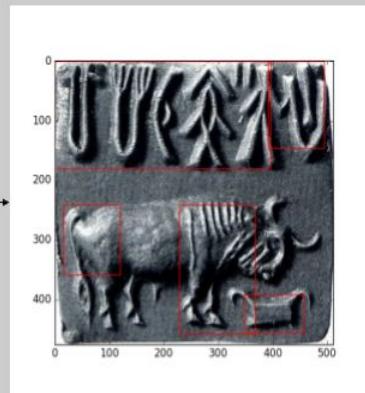
Indus Seal Scan



Extract Seal



Selective Search



Region Grouping

# EXTRACT SEAL

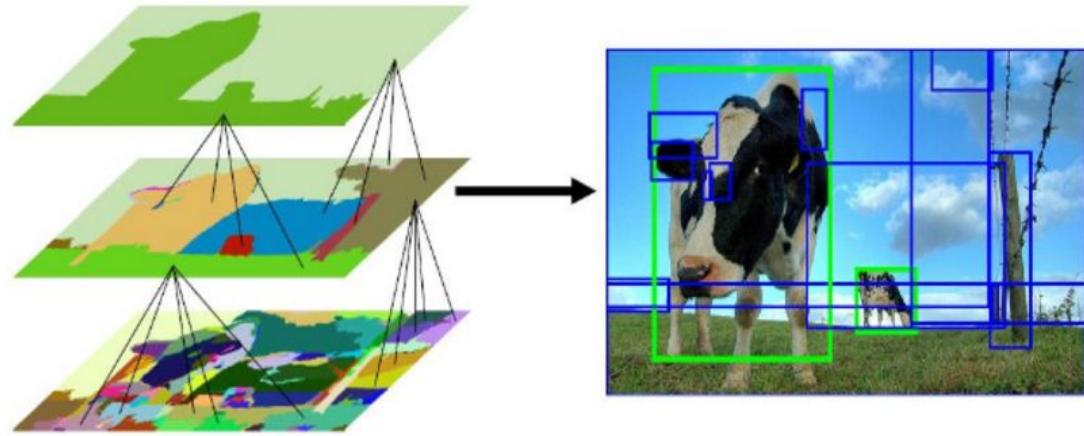
Removes the irrelevant background information from the input artifact image.

- Steps
  - Grayscaled and smoothed 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

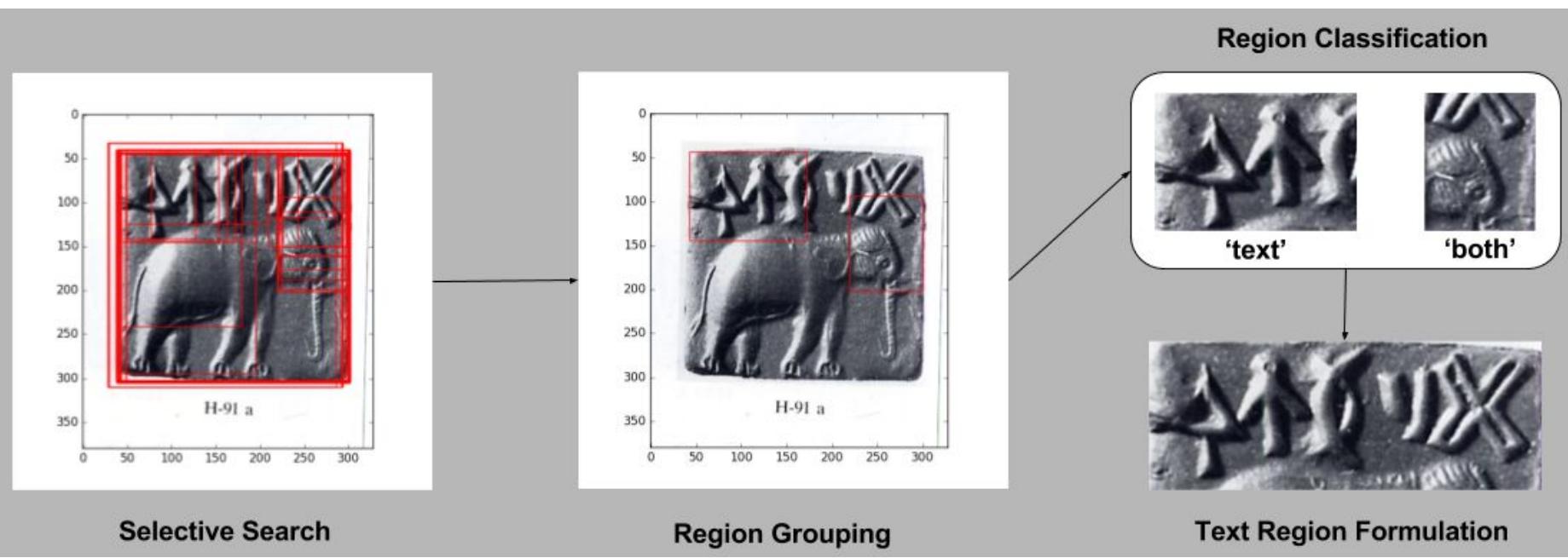
Improvises the quality of the regions proposed by Selective Search (four-level region grouping and filtering hierarchy)

- 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

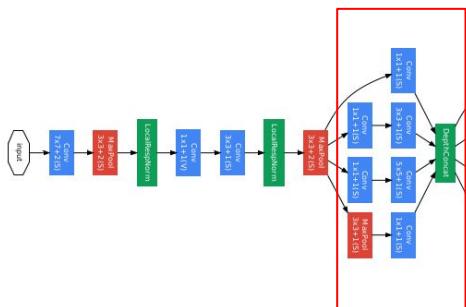
Produces exact text-only regions by eliminating the non-graphemic parts off of the region proposals in hand.



# REGION CLASSIFICATION

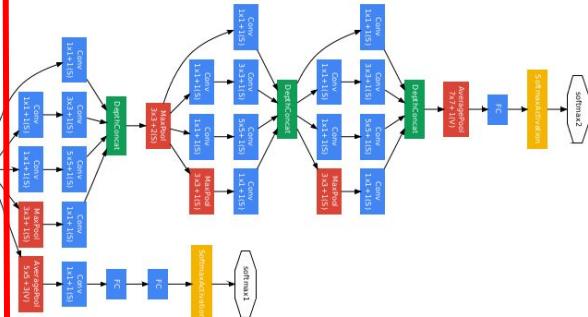
All the regions are classified into types: “Text”, “No-Text” or “Both”.

ImageNet weights - Frozen (Transfer Learnt)



Inception Module

ImageNet weights - Initialized



Learning Rate - 2X (Fine-tuned)

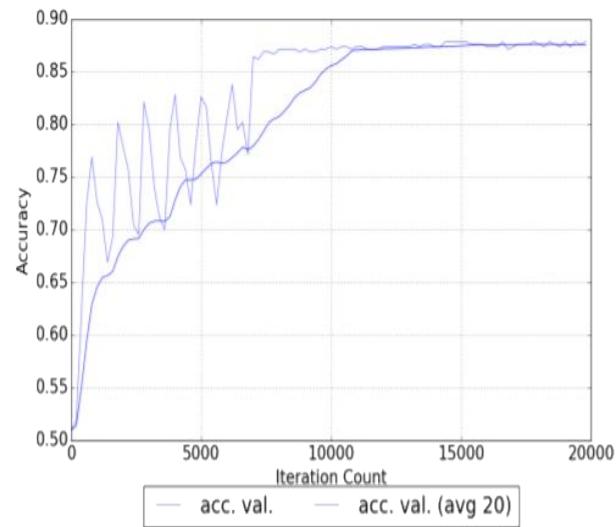
# THE REGION CLASSIFICATION CNN'S RESULTS

## Top-1 Accuracy Scores

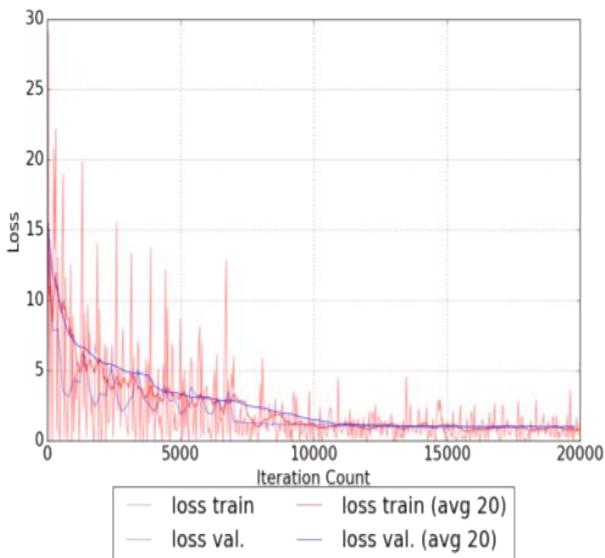
GoogLeNet's Levels	Top1 Accuracy Scores
Level 1 (1/3 <sup>rd</sup> network depth)	87.14%
Level 2 (2/3 <sup>rd</sup> network depth)	87.86%
Level 3 (full network depth)	89.30%

		
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

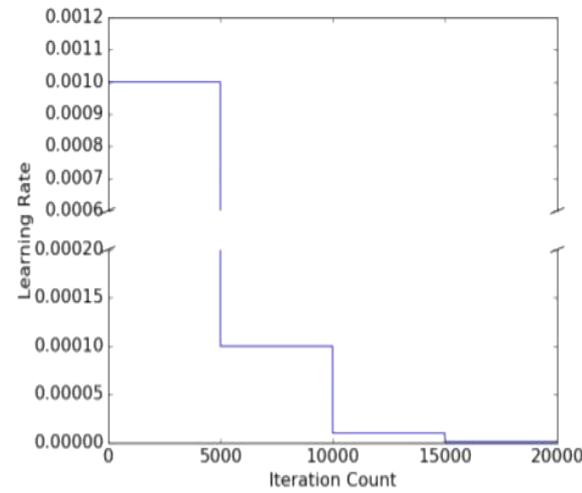
# THE REGION CLASSIFICATION CNN'S GRAPHS



(a)



(b)



(c)

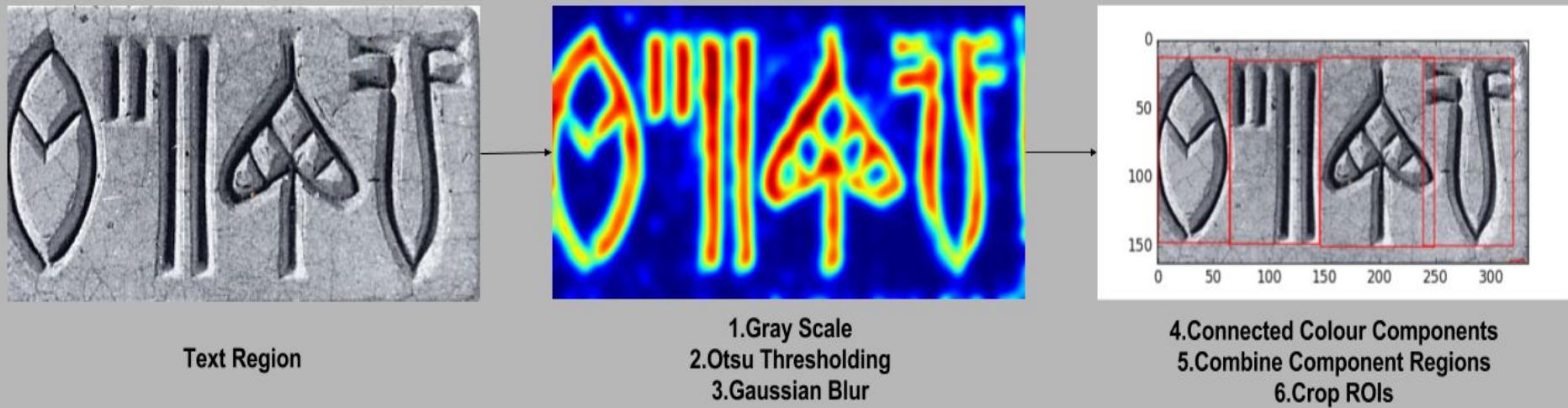
# TEXT REGION FORMULATION

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.



# SYMBOL SEGMENTATION - THE ALGORITHM

---

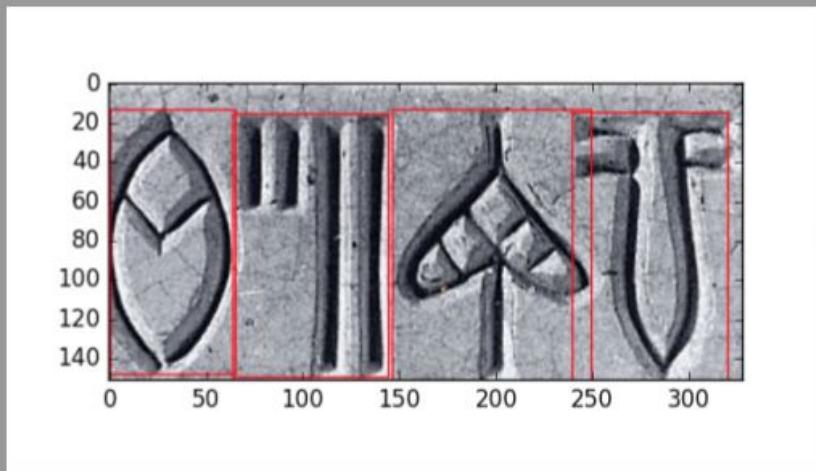
## Algorithm 1 Symbol Segmentation Algorithm

```
1: procedure SEGMENT-SYMBOL(Image I)
2:     Gray_Image = Gray_Scale(I)
3:     Thresholded_I = Otsu_Thresholding(Gray_Image)
4:     Smoothened_Image = Gaussian_Blr(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
```

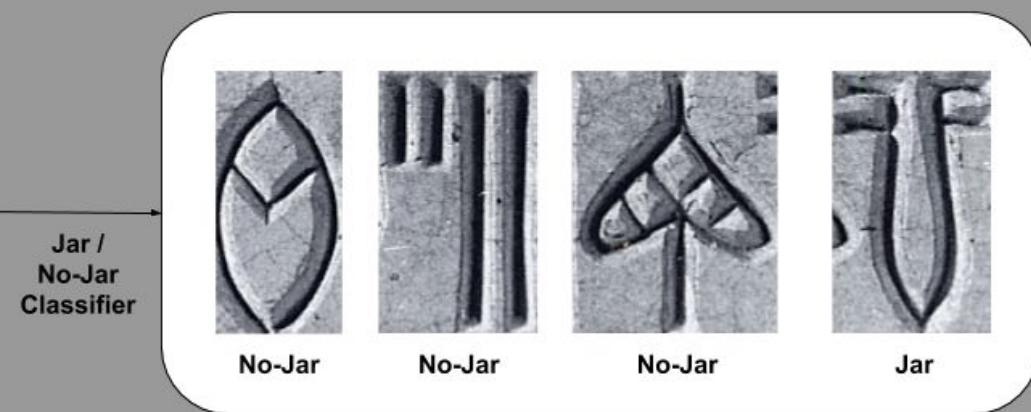
---

# SYMBOL IDENTIFICATION

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



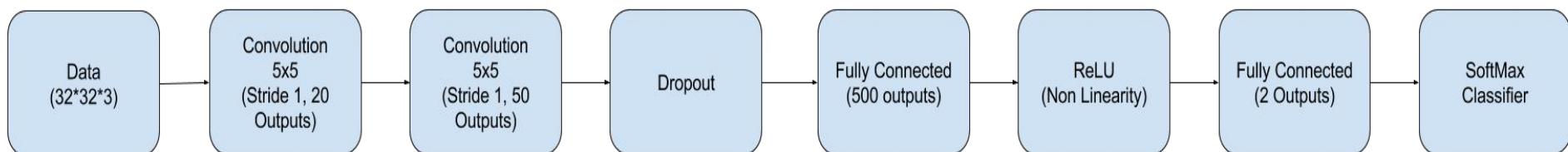
Symbol-wise Segmented Text Region



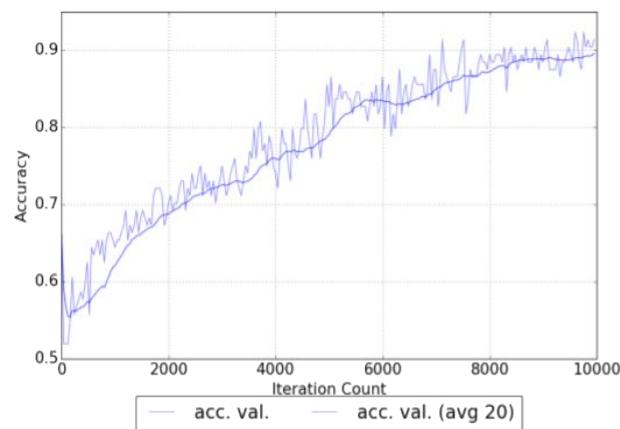
Classified Symbols

# SYMBOL IDENTIFICATION - CNN ARCHITECTURE

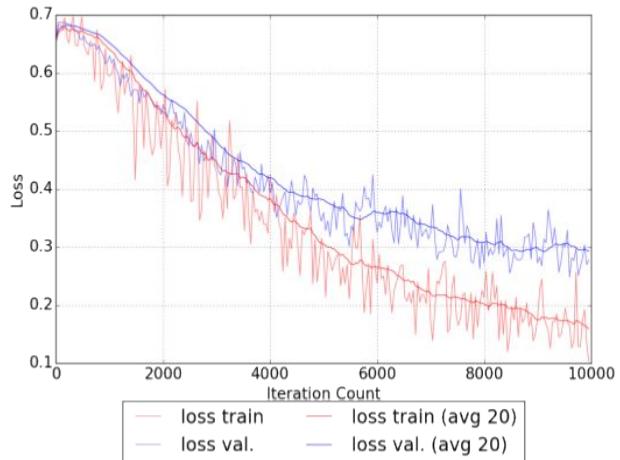
- Detects the presence or absence of the “Jar” sign – Binary classifier
- No transfer learning
- New architecture trained from scratch
- Accuracy score of **92.07%** when evaluated over the validation set of the “Jar-NoJar Dataset”



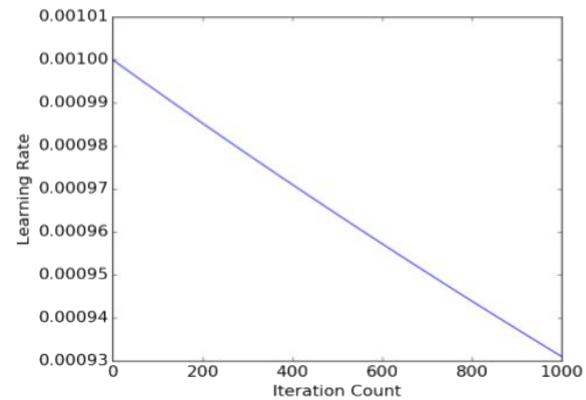
# SYMBOL IDENTIFICATION CNN'S GRAPHS



(a)

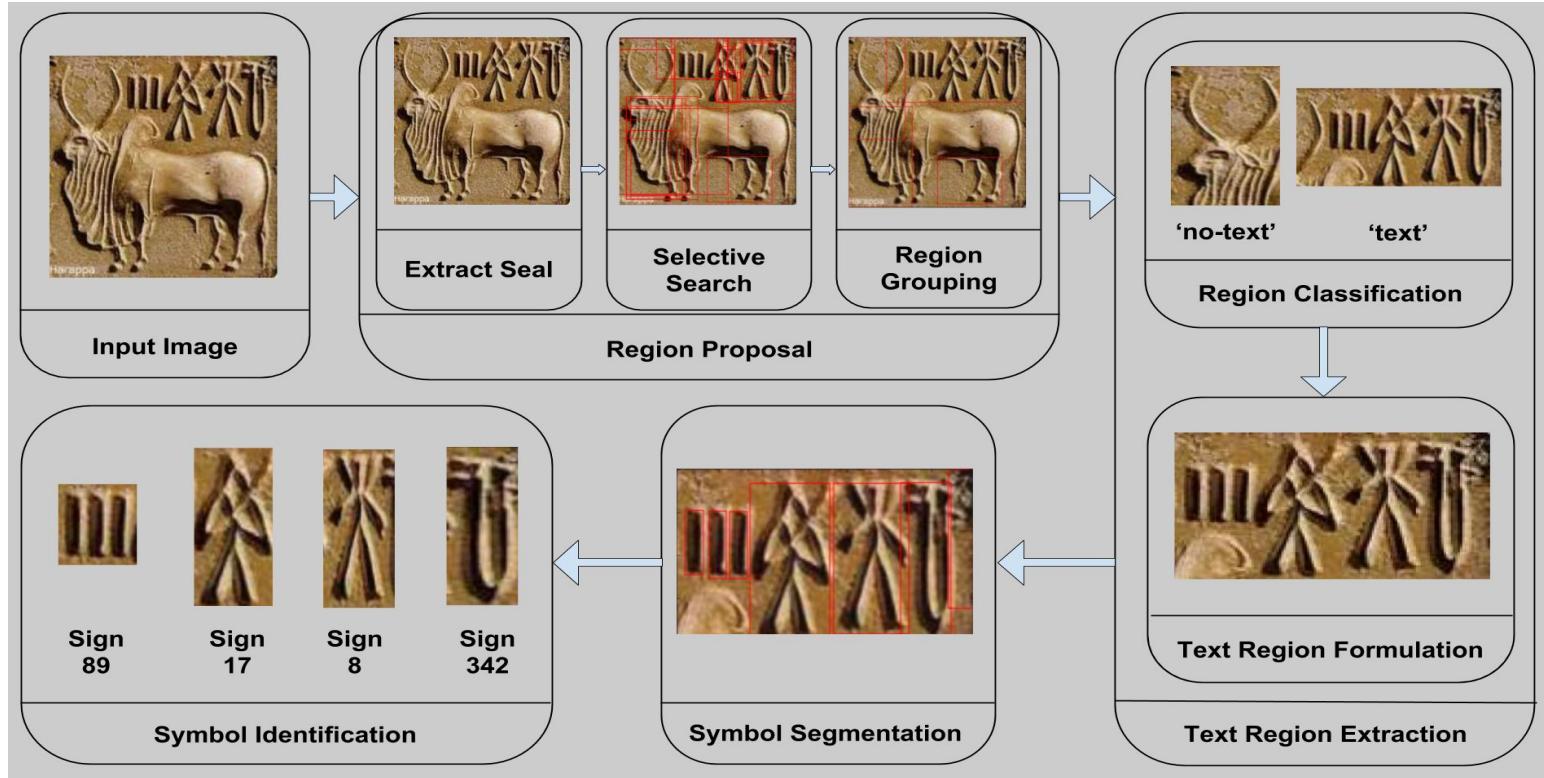


(b)



(c)

# AN EXAMPLE FLOW

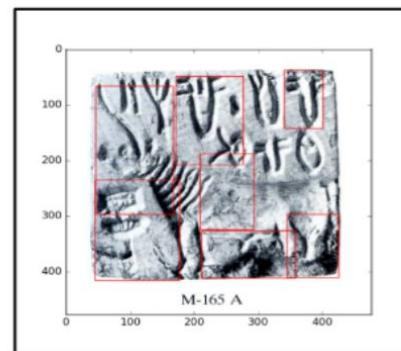
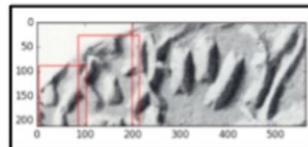
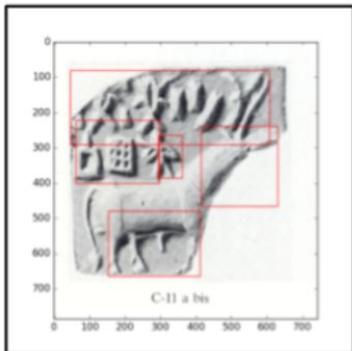


# EVALUATING THE PIPELINE

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			
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

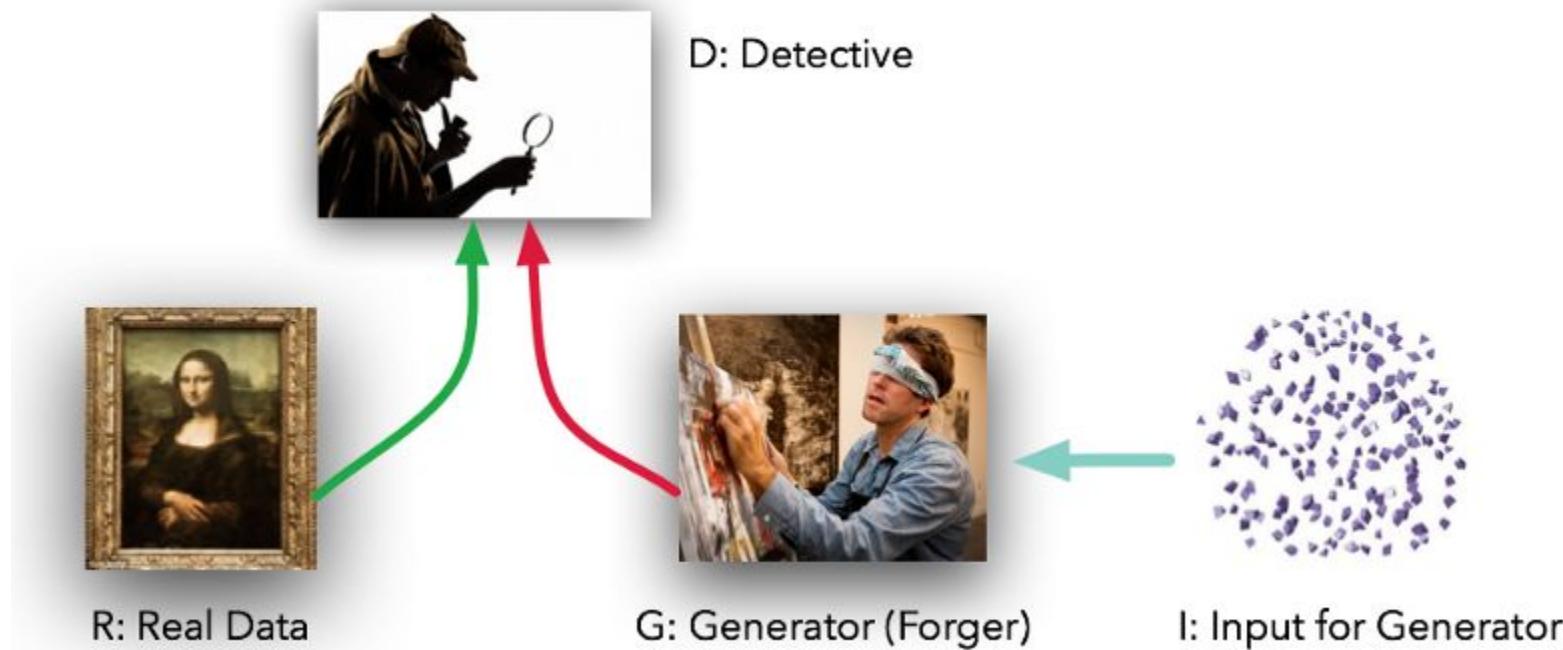
- Non ML based techniques used in Region Proposal and Symbol Segmentation stages
  - Leads to some parts of the text region and symbols getting missed
- Extremely damaged seals can never be fully read
- Very complex grapheme arrangement structure
- Limited to identify only few frequent symbols of the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

# FUN TASK: GANS



# THE BUZZ !



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

By Mallory Locklear | July 25, 2017, 9:00am EST

Meanwhile, Ronjooy 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.

IN FOCUS Eurovision Life Indigenous Sexuality The Playlist podcast  
SBS HOME ON DEMAND GUIDE PROGRAMS RADIO SHOP NEWS CYCLING FOOTBALL MOVIES  
Tamil home News & features Tamil Team Contact us

5 MAR 2017 - 8:26PM  
Translation in English தமிழ்  
**An app to Decipher ancient symbols**  
PODCAST tamil.170305.641975.mp3  
00:00 08:15

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 Kulasegaran Sanchayan about it.

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

Utejomayam  
@timegroup.com

**A**s the Egyptian civilisation flourished, Indian seal engravers documented the rise and fall of one of man's greatest civilisations in the Indus Valley, c.3,500BC-1,300BC. Another great civilisation arose in the Indus Valley in the northwest of the Indian subcontinent.

Much less is known about the Indus Valley Civilisation than the Egyptian, however, about its development, governance, activities, discoveries and daily life – because historians have been relatively short on information and yet to fully interpret the script of these ancient pieces of art.

To the common man, however, the limited corpus of hieroglyphs and other symbols that have been found in artifacts from the Indus Valley bear an uncanny resemblance to those found along River Nile. Scientists believe the Indus Valley had Dravidian languages and an early form of Sanskrit but its meaning remains a mystery.

Here's a discovery, however, that could help change that.

Artificial intelligence involves a computer system that drives cars, which mimics the functions of the human brain, may now aid researchers to develop a more accurate database for Indus script that could eventually help decipher the texts. Scientists are also working on a mobile application for the software.

The technology will allow archaeologists and amateur history buffs alike to, say, capture images of seals on pottery and share them 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. The app will filter the text from the image and identify the sequence and sense of individual characters in an existing database. If it is a known symbol, the app will display a number representing each character in the texts in

### DAWN OF TIME

The computer application can be used to identify elements belonging to the Indus script

> The image of the seal is scanned on all sides that it has to have the Indus script symbols, depictions of animals like bulls and unicorns, and deities

> A modern take on history

> The Indus Valley Civilisation is one of the oldest known of the ancient civilisations. It existed from 2500 BC to 1300 BC

> Discovered in the 1850s, it was spread over 1.2 million sq km covering parts of modern-day Pakistan, Afghanistan and India

> It included around 1,000 well-planned villages, towns and cities, with

> The image is further classified into 'text', 'no-text', and 'seal'. The seal consists of only the Indus script non-graphemic elements.

> A customised algorithm segments out individual units or letters to identify it

> The identified units or letters are classified into one of 417 classes, the known names of Indus graphemes, according to scholar Mahadevan

> AS recently excavated structures resembling those at Harappa from Keeladi village, Tamil Nadu. They are believed to be 2,000 years old



> Two prominent cities, Mohenjo-daro and Harappa, and Mohenjo-daro  
> Seals were used for ritualistic, commercial and religious purposes  
> Copper, bronze, pottery and terracotta toys, beads (carnelian) and (steatite) seals were engraved with animal figures  
> Though the origin and decline of the civilization remain a mystery, there are theories that Greeks, defeated after the invasion of Aryans could be among reasons

the database if not, it will include the symbol in the database.

The output will be a string of graphemes (characters) and a corresponding number. This will be for data accumulation and analysis. Symbols. Automating corpora preparation will speed up research to decipher the script.

A professor from Institute of Mathematical Sciences (IMSc), Chennai, and engineer have developed the app. Satish Palaniappan, an SSN College of Engineering graduate who worked on the app, said the sequence of numbers may help in the search for similar sequences in other texts. Palaniappan and Adhikari used 'deep learning' to develop the technology. "Deep learning is based

between regions of the Indus Valley.

The app is crucial to make big leaps in epigraphic research. "A researcher has to know the history and sequence of symbols," Palaniappan said. "It takes years to compile texts from artefacts and put them in form. The computer will understand. We wanted to bridge that gap."

## Chennai team taps AI to read Indus Script

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

SHUBASHREE DESIKAN

The Indus script has long challenged epigraphists because of the difficulty in reading and classifying text and symbols on artefacts. Now, a computer-based team of scientists has built a programme which eases the process.

Ronjooy Adhikari of the Institute of Mathematical Sciences and Satish Palaniappan, who is at Sri Sivasubramanian Engineering, have developed a "deep learning" algorithm that can read the Indus script from images of artefacts such as seals or pieces that contain Indus writing.

Scanning the image, the algorithm recognises the region of the image that contains the script, breaks it up into individual units and then feeds it into a database for the smallest unit of the script and finally identifies these using a deep learning algorithm. In linguistics the term corpus is used to describe a large collection of texts while, among other things, are used to train computers to analyse lyrics of languages.

The process consists under a class of artificial intelligence called "deep neural networks." "These have been a major part of the game-changing developments in fields of driving cars and Go-playing bot that surpass human performance," says Mahadevan. The deep neural network mimics the working of the mammalian visual cortex, known as convolutional neural network (CNN), which divides the field into overlapping regions. The features found in each region are hierarchical, so that they can be used to build a composite understanding of the whole picture.

**Indus script**  
The Indus Valley script is much older than the Brahmi and Tamil-Brahmi scripts. However, unlike the latter two, it has not yet been de-

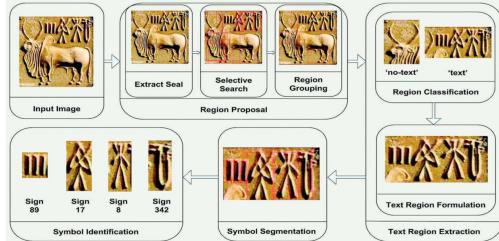
ciphered because a bilingual text has yet been found.

A bilingual text has in many other cases aided archaeologists in understanding ancient scripts, for example, the Rosetta stone, which was found in the eighteenth century carries inscriptions of a decree issued in 190 BCE, in three parts, the first two in an ancient Egyptian hieroglyphic and the Demotic script, while the bottom is in Greek.

Since the day the Rosetta stone was provided to the British Museum, the work on deciphering the Indus script has been progressing. For the past 40 years, the Rosetta stone's Indus script remains undeciphered today.

It is a major effort to even build a standard corpus to the language and decode the writing on existing artefacts and map them to this standard corpus. The most widely accepted theory is that the Indus script was brought together by the efforts of Iravatham Mahadevan, noted Indus scriptologist, from the 3,700 texts or 417 unique characters he compiled about the relevance of this work, 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."

He has asked about the relevance of this work, 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."



**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

It is a major effort to even build a standard corpus to the language and decode the writing on existing artefacts and map them to this standard corpus. The most widely accepted theory is that the Indus script was brought together by the efforts of Iravatham Mahadevan, noted Indus scriptologist, from the 3,700 texts or 417 unique characters he compiled about the relevance of this work, 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."

**Dr Mahadevan says, "It 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

