

Summer Internship

Centre for Development of Advanced Computing
CDAC, Pune

A presentation by
Kimya Gandhi
08625005

IDC, IIT Bombay

Objective of project

- To study font design in Indian scripts
- To study Indian scripts and their diversity
- Get acquainted with softwares and standards
- Learn the process of font design; the research and methodology required.

Introduction to CDAC & GIST

- Centre for Development of Advanced Computing (CDAC) has been a pioneer in developing and proliferating use of Indian languages on the computer.
- Their Graphic and Intelligence based Script Technology (GIST) has been a a bridge between Indian languages and IT
- They are in tune with technologies worldwide
- It is well known for its research and digitization of Indian languages and making it available to the masses.
- The division brags a collection of fonts in the various scripts of India, Devanagari, Gurmukhi, Bengali, Tamil, Telegu, Malyalama, Gujrati, Sinhalese to name a few.



Project Outline

- Project: Font Design
- Script: Devanagari
- Application: Optical Character Recognition (OCR) systems
- Understanding OCR technology and its application
- Study of devanagari script and nature of its characters
- To create a font that would be specifically give accurate results for character recognition.

Digitization and standards

- Need for standards: because one font is used over diverse applications and softwares in different places.
- The need for digitization of Indian scripts: To increase their reach and prevent them from being obsolete
- The use of one script for writing different languages
- UNICODE, ASCII, ISCII
- Opentype font formats

Hex		0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
Hex	Dec.	0	16	32	48	64	80	96	112	128	144	160	176	192	208	224	240
0	0	NUL	DLE	SP	0	@	P	.	p								
1	1	SOH	DC1	!	1	A	Q	a	q			ॠ	औ	ष	ॡ	ॢ	EXT
2	2	STX	DC2	"	2	B	R	b	r			ॡ	ऑ	ॢ	ॣ	।	०
3	3	ETX	DC3	#	3	C	S	c	s			ॢ	क	ॣ	।	॥	१
4	4	EOT	DC4	\$	4	D	T	d	t			ॣ	ख	।	॥	०	२
5	5	ENQ	NAK	%	5	E	U	e	u			।	ग	॥	०	०	३
6	6	ACK	SYN	&	6	F	V	f	v			॥	घ	०	०	०	४
7	7	BEL	ETB	'	7	G	W	g	w			०	ज	०	०	०	५
8	8	BS	CAN	(8	H	X	h	x			०	घ	०	०	०	६
9	9	HT	EM)	9	I	Y	i	y			०	ञ	०	INV	०	७
A	10	LF	SUB	*	:	J	Z	j	z			०	ट	०	०	०	८
B	11	VT	ESC	+	:	K	[k	{			०	ड	०	०	०	९
C	12	FF	FS	.	<	L	\	l				०	ढ	०	०	०	
D	13	CR	GS	-	=	M]	m	}			०	ॢ	०	०	०	
E	14	SO	RS	.	>	N	^	n	~			०	ॣ	०	०	०	
F	15	SI	US	/	?	O	_	o	DEL			०	।	०	ATR	०	

8-bit Code Table of the Latin and Indian Script Alphabet

Optical Character Recognition (OCR)

- Understanding Optical Character Recognition (OCR): it is the mechanical or electronic translation of images of handwritten, typewritten or printed text into machine-editable text.
- OCR for Latin script: a largely solved problem
- Applications of OCR: credit cards, postal services, data entry, preserving and digitizing old documents or legal documents
- Reasons for using OCR
 - Reduce data entry errors
 - Consolidate data entry
 - Handle peak loads
 - Human readable
 - Many printing techniques
 - Scanning corrections



Cheque number printed in OCR readable font

Devanagari script

- An overview of the script: written from left to right, no distinct cases, *shirorekha* or the top line is a distinct feature, based on phonetic structure
- Devanagari in manuscripts: a brief study of the precedents of the present day used devanagari.
- Analysing different writing systems
- Nature of the script and alphabet: consonants, conjuncts, vowels, vowel signs, diacritic marks, punctuations.
- Anatomy of the devanagari letterforms

Consonant Signs with Following ă					
Velars	क	ख	ग	घ	ङ
	ka	kha	ga	gha	ṅa
Palatals	च	छ	ज	झ	ञ
	ca	cha	ja	jha	ña
Linguals	ट	ठ	ड	ढ	ण
	ṭa	ṭha	ḍa	ḍha	ṇa
Dentals	त	थ	द	ध	न
	ta	tha	da	dha	na
Labials	प	फ	ब	भ	म
	pa	pha	ba	bha	ma
Semivowels	य	र	ल	व	
	ya	ra	la	va	
Sibilants	श	ष	स		
	śa	ṣa	sa		
Aspirate	ह				
	ha				

Devanagari font design

- Font design in devanagari is a complex process: no complete set.
- The non linear and complex nature of the script
- The difference in physical form of the same letter used in different languages
- Inappropriate placement of vowel signs (*matras*), consonant clusters form complex conjuncts and ligatures.

मुद्राक्षरलेखनकला

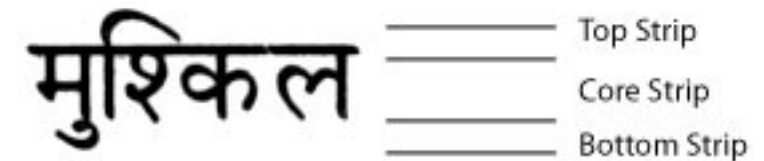
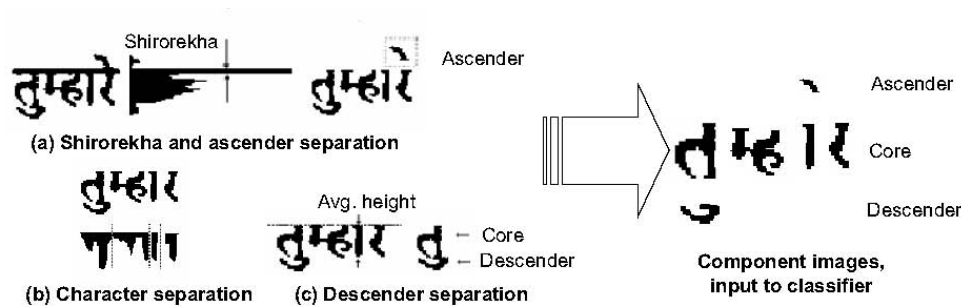
Non linear nature of script

अल्पित नूपुर

Inappropriate placement and spacing of matras

OCR and Devanagari script

- Chitrakan: OCR software for Indian languages
- Process of character recognition: feature extraction, segmentation
- Problems with devanagari character recognition
 - recognition difficult in smaller point sizes
 - ambiguity in characters with similar structures
 - non linear nature of script
 - shirorekha breaks after a word as also in characters like अ, श, झ
 - ligatures and complex conjuncts



Segmentation for recognition; Chitrakan

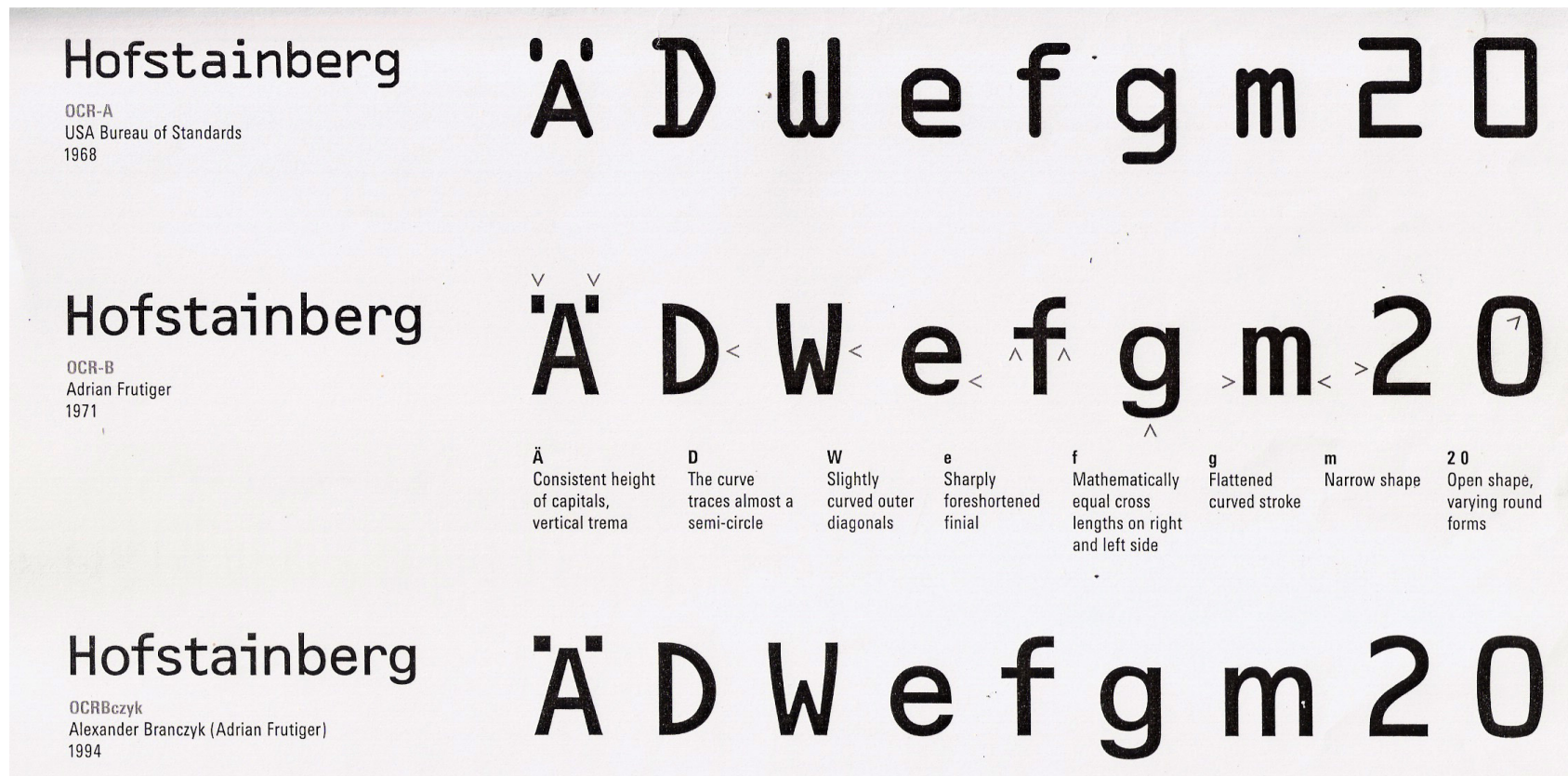
Segmentation driven OCR

Character design

- The study of problems faced by recognition of devanagari script led to setting parameters for font design
 - characters that could be scanned with accuracy and speed
 - simple and open counters
 - readable in small point sizes
 - machine and human readable
- The stages involved in the design of letterforms included following stages
 - case studies; OCR specific fonts
 - visual references; understanding proportions
 - study of existing fonts
 - explorations for construction of letterforms
 - character set design

Case studies

- OCR-A, OCR-B
- Understanding methodology and logic that went behind designing the font



Above: Comparison between OCR-A and OCR-B

Visual references: Letterpress

Letterpress printed ads and book covers

वरप्रार्थना : पञ्चांगदान

महाराष्ट्र राज्य शासकीय प्रकाशन

संस्कृतित श्रीश्रीश्री विकास
कान्तिपत्रिका अखिलीक वित्तवृत्त वकी विज्ञान

लोकशक्ति

मुंबई राज्याचे शिलकी अंदाजपत्रक इजिप्तमध्ये राण्यकांति
योजना कार्यावर भरपूर खर्च
अचमकी ही, जीवसाजोडिक्मेमाडी वया इति
मिहना शोधे मापय मिडोपनामाडी प्रामांनिडी

गणतन्त्र
नर्गावचा राजीनामा

विधायने वार करन बापकोचा
जमानप यून

बी. एम्. गोखले
अँड सन्स

ऑप्टिशियन्स -
कॉन्टॅक्ट लेन्स
क्लिनिक



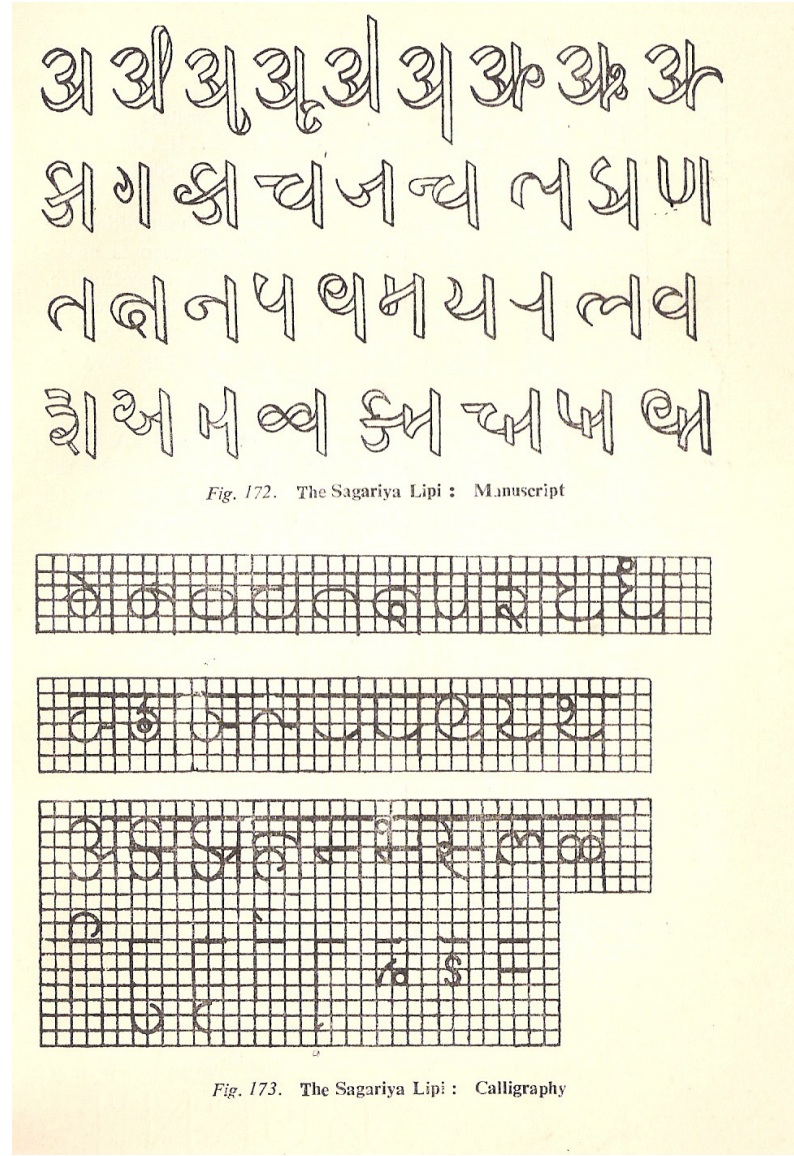
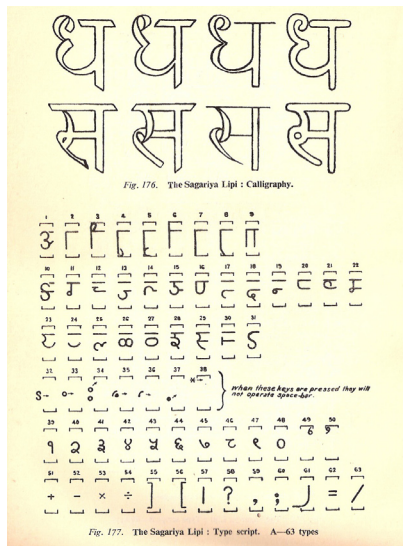
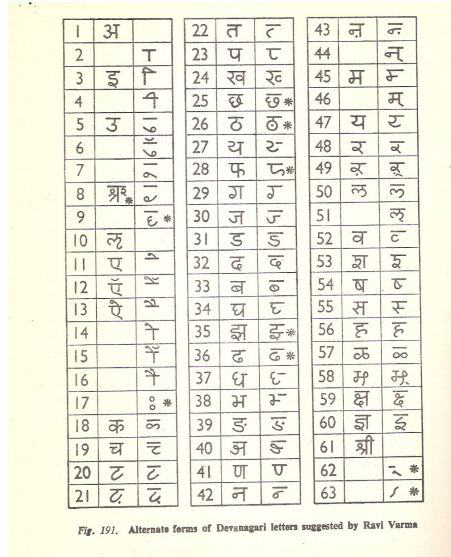
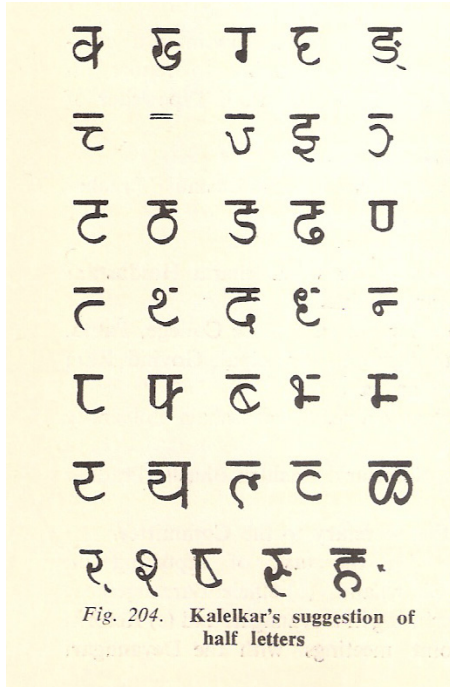
चांदिका महाड, शिवाजी पार्क जवळ, ले. ज. रोड, दादर, मुंबई-४०००२८

भावी काल

ज्ञानदेवांचे
शास्त्रप्रमाण
तर्कतीर्थ लक्ष्मणशास्त्री जोशी

हवेली तालुका कांदा उत्पादकांचा
सहकारी खरेदी-विक्री संघ लि.
छत्रपती शिवाजी मार्केट यार्ड, गुलटेकडी, पुणे-३७.

Visual references: Writing systems



Linear writing systems

Other writing systems and proportions.

Visual references: Typewritten text

यातील पक्षाकरात असेही मान्य व सहमत करण्यात येते की,
कोणाही पक्षाकाराने निवृत्ती स्विकारल्यास वा सदरील
पक्षाकार हा निधान पावल्यास सदर भागीदारी व्यवसाय
हा बरखास्त करण्यात येणार नाही, उर्वरित पक्षाकार
भागीदारी व्यवसाय, योग्य प्रकारे चालविण्यात येणार आहे
व त्यानुसार त्या त्या योग्य असे वाली / वारसास त्या

Above; Devanagari text typewritten page section & right; the
entire page

Observations:

Quality of printed page; no grey mass

Spreading and blotting of ink in complex characters or
ligatures

Matras and their placement, spacing

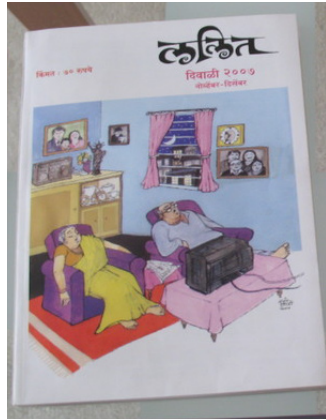
: ६ :

आणि नफा तोट्यापत्राके बनविण्यात येणार आहेत व त्यावर
पक्षाकार आपआपल्या स्वाक्षा-याही करणार आहेत व
त्यानुसार काही चुका, दुस-या काही बाबी असल्यास त्या त्या
आपसात मान्य व सहमतीने हाताळण्यात येणार आहेत. व
त्यानुसार पक्षाकार आपआपल्या स्वाक्षा-याही त्यावर करतील
वा सुधारीत बाबीचे नुसार लिखित कागद पत्राकावर स्वाक्षा-या
करून सहमतीने ते ते व्यवहार हाताळण्यात येणार आहेत.

१३. यातील पक्षाकारांत असेही मान्य व सहमत करण्यात येते की सदरील
भागीदारीचे अन्वये नावलौकीकता तसेच कोटा अणिकार
लायसन्स, दुरध्वनी आणि इतर टेलिफोन कनेक्शन व त्यानुसार
नपघाचे व्दारा काही बाबी असतील त्यानुसार अधिकार राहिल.

१४. यातील पक्षाकारात असेही मान्य व सहमत करण्यात येत आहे
की, जर कोणाही पक्षाकारात निवृत्ती स्विकारण्याचा
असल्यास त्या पक्षाकाराने उर्वरित पक्षाकारात किमान सहा
महिने आधी लेखा निवेदन सादर करण्याचे आहे व
त्यानुसार त्या त्या पक्षाकाराने निवृत्ती स्विकारण्याची आहे.

Visual references: Book covers, learning aids



Visual references: Existing fonts

यह सीडॅक की योगेश मुद्रलिपि है।
निम्नलिखित
बीजाङ्कुरित
काँङ्करर
काकचेष्टा बकोध्यानम्
ष्ठर्योः |

DV-OT Yogesh designed by GIST team, CDAC

अआइईउऊऋऌएँऐएरेऑऋऌऌऌऌऌ

Mangal™ (Devanagari)

खगघङ्चछजझञटठडढ
णतथदधनपफबभमयर
लळवशषसहक्षज्ञक्रगघङ्

Mangal designed by Prof. R. K. Joshi

Character design

- The process of designing characters included:
 - character listing; *varnamala*
 - categorizing characters according to various parameters
 - openness, number of strokes, joineries with shirorekha
 - geometric construction of form
- The study of OCR, devanagari and visual references led to setting parameters for design of the font which are:
 - geometric construction
 - uniform weight and width of characters
 - one pixel gap between characters
 - simple joineries
 - open counters

Character segregation

number of strokes

अ आ इ ई उ ऊ ए ऐ ओ औ अं अः

ण ध ट थ त्र घ छ ष क्ष ख भ र त य प स द फ ग ह ज क ल झ श

च व ब न म ड त् श्र ज्ञ रु रू ढ ऋ

one stroke

र ट द ठ ढ ड इ उ

two stroke

ण ग ह ज भ म य प त न व ल च थ छ ए ऊ ई ळ श ज्ञ

three stroke

अ स क ख फ झ त्र ष क्ष ब ऋ

Character segregation

points of contact with *shirorekha*

one point

र ट द ठ ढ ड इ उ ह ज त न व ल ई ऊ ऌ च ब क त्र ज्ञ ऋ

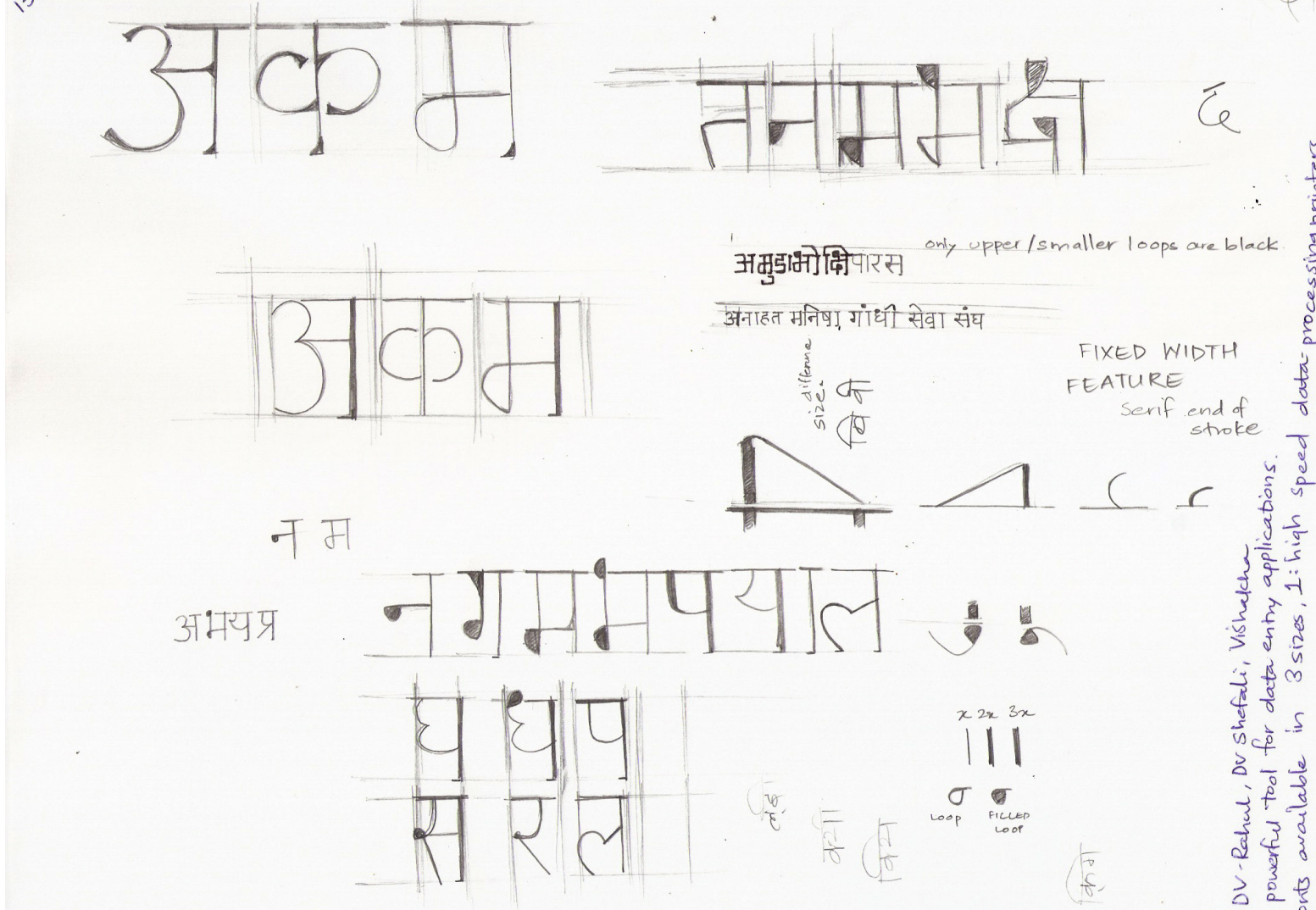
two point

ग म य प छ ए स ख फ झ ष

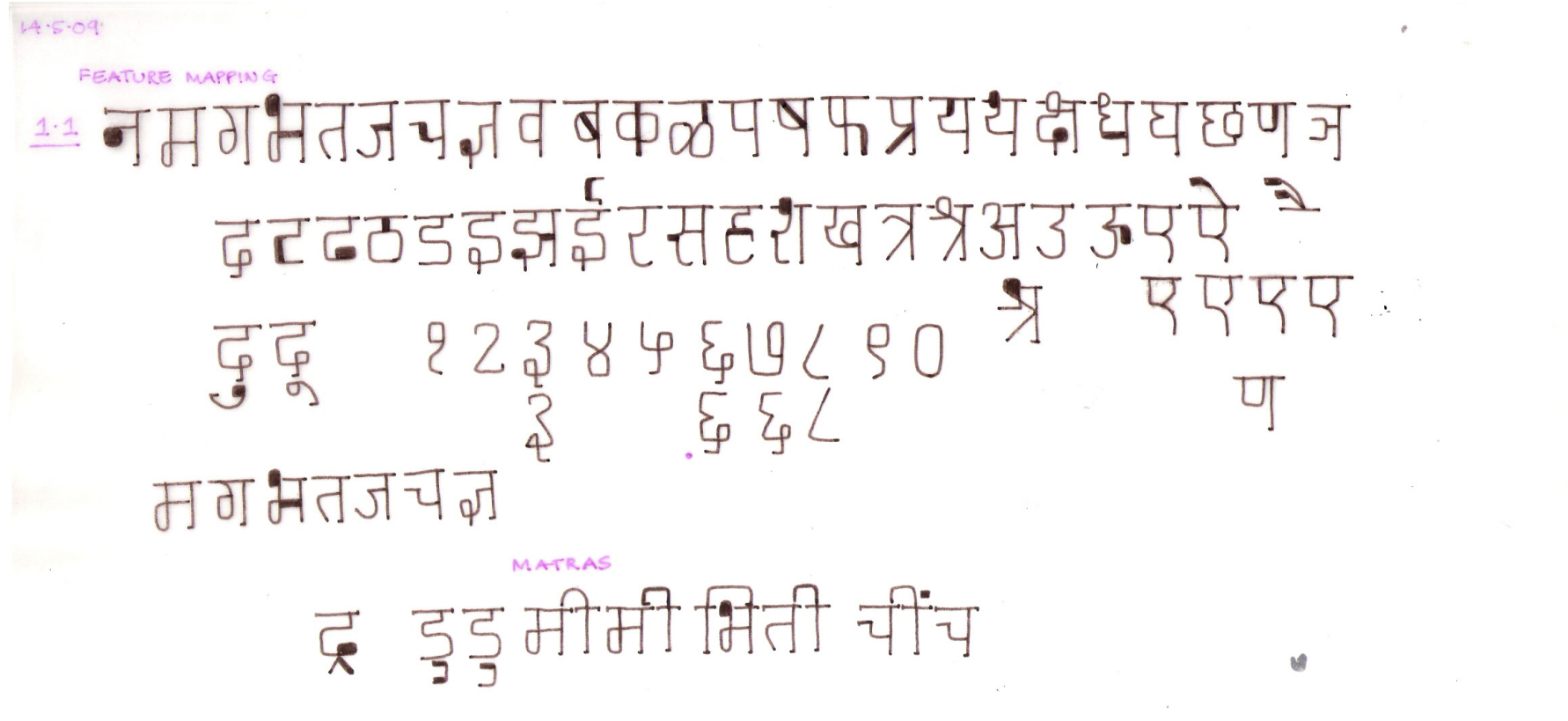
three point

अ क्ष ण भ थ श

Character design; form exploration



Character design; geometric construction



Characters based on lines and slight curves.

Form inspired from the construction of OCR-A

- Letterforms lack the devanagari character

Character design; geometric construction

1.4

SQUARES
AND
LINES

क ट ठ ड इ झ र स ह श य त्र श्र

अ उ ए ऐ

म ग म त ज च ज्ञ व ब क ङ प ष फ प्र य थ क्ष घ घ छ ण ञ

दु दु मी मिती चींच ३ ३ MATRAS दृ द्र द् द्र

१ २ ३ ४ ५ ६ ७ ८ ९ ०

१ २ ३ ४ ५ ६ ७ ८ ९ ०

श

1.5

क ट ठ ड इ झ र स ह श य त्र श्र

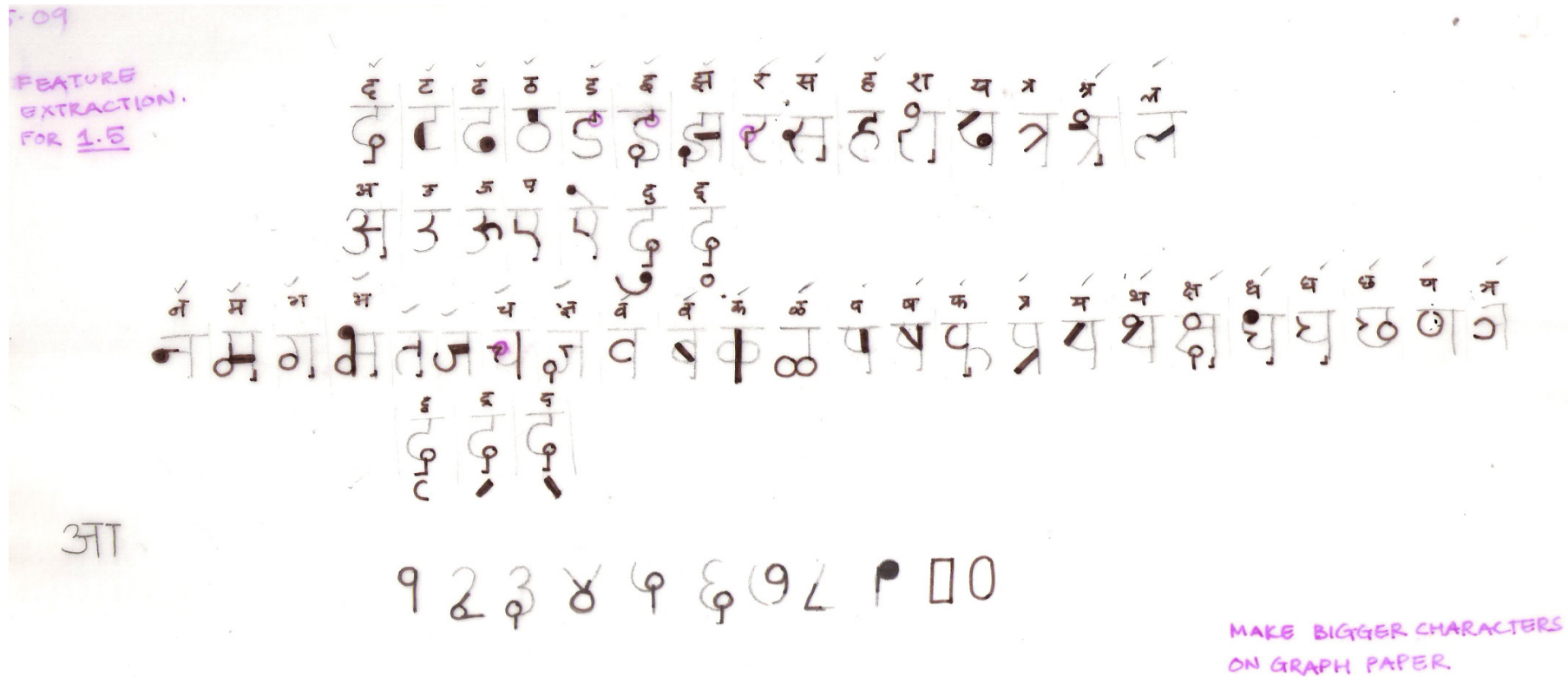
अ उ ऊ ए ऐ दु द्र

न म ग म त ज च ज्ञ व ब क ङ प ष फ प्र य थ क्ष घ घ छ ण ञ

The circle makes characters more open, readable and restores its natural form

Feature extraction

Feature extraction; highlighting features that distinguish one character from the other.

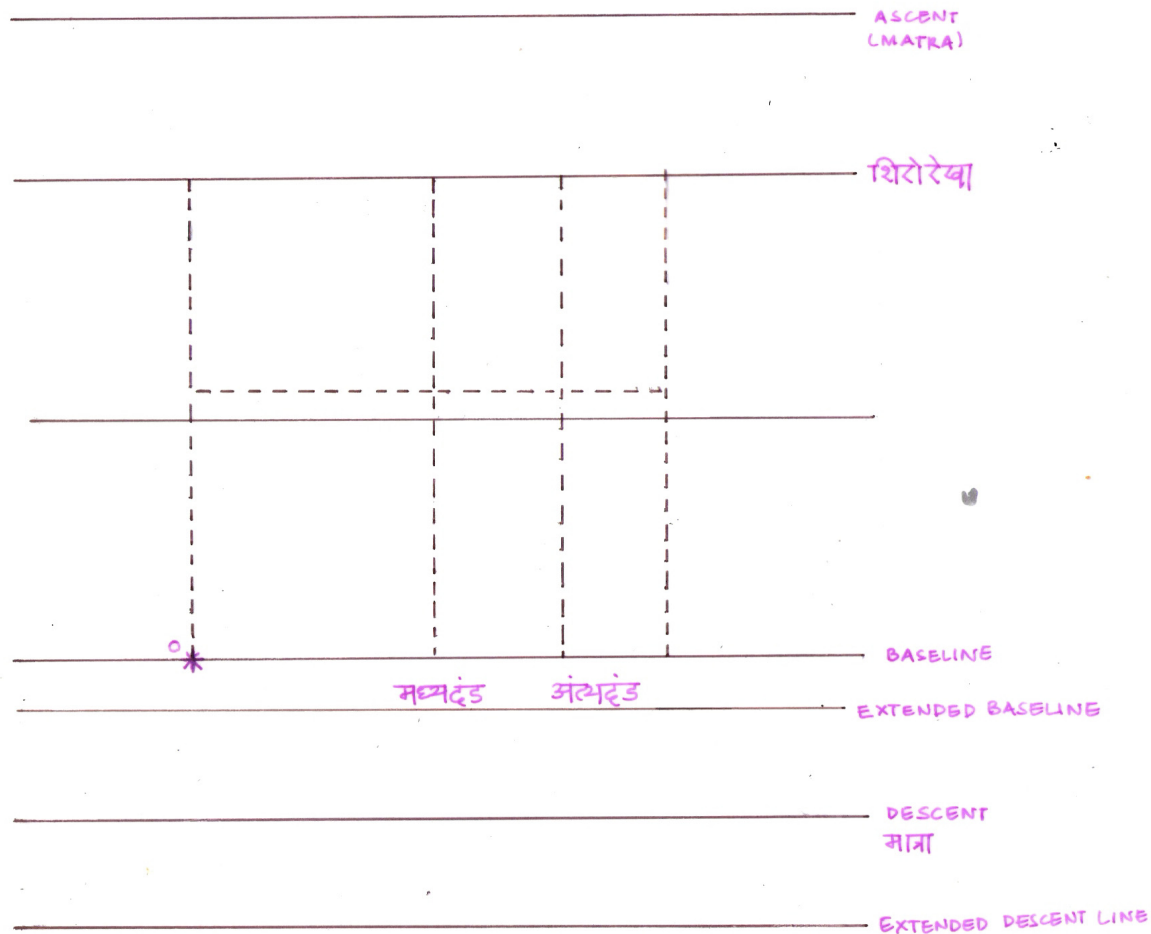


This type of a circle and line construction was finalised and letterforms were constructed on graph sheets with proper measurements and a grid.

The Grid; circle and line

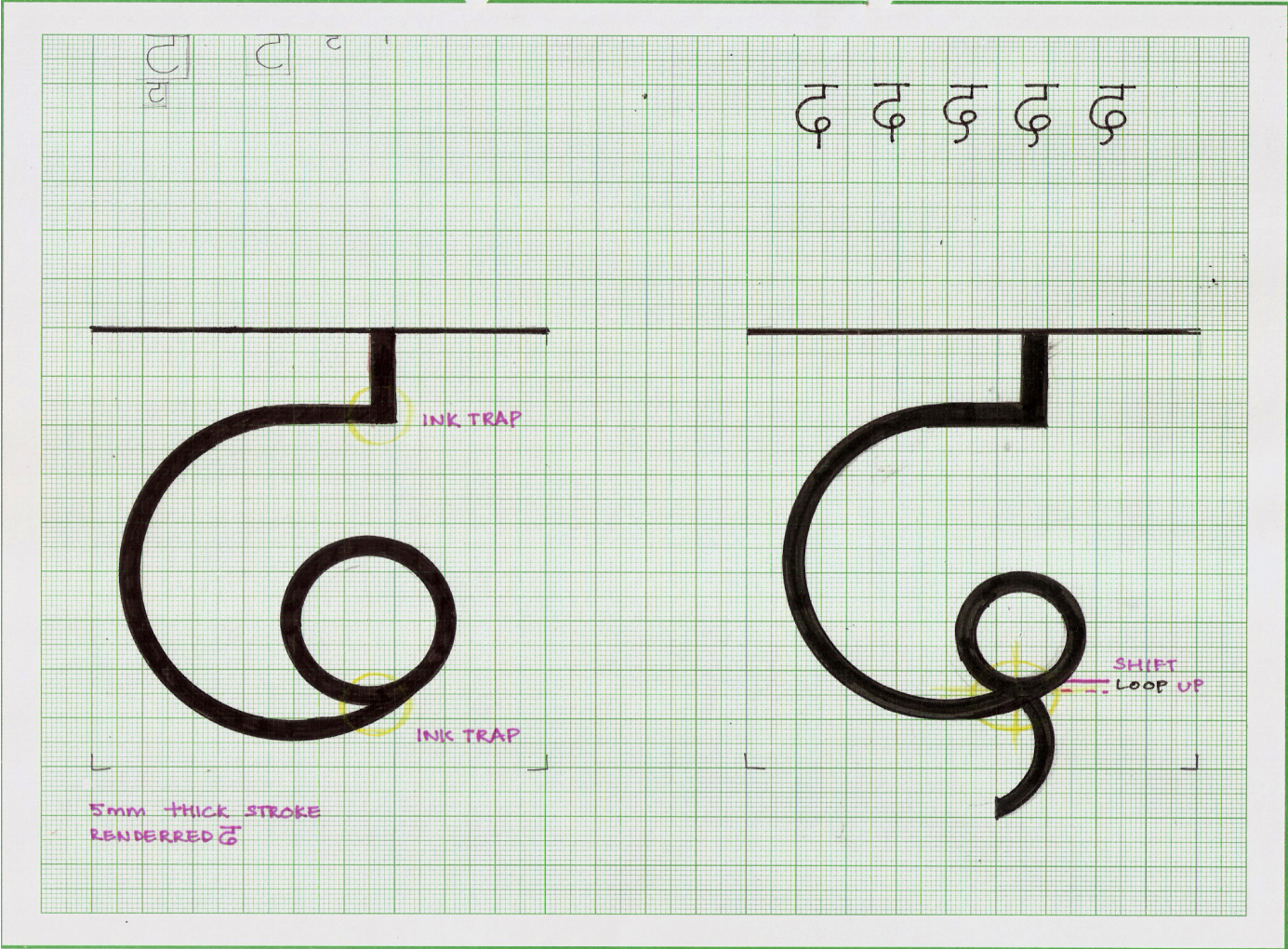
The grid defines framework of the character set

The height to width ratio is close to one

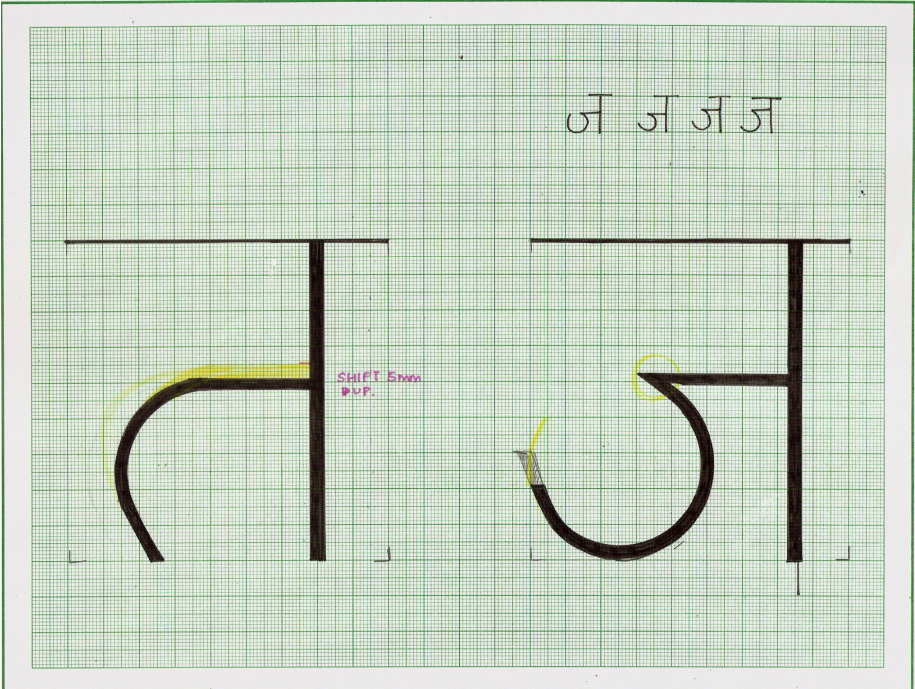
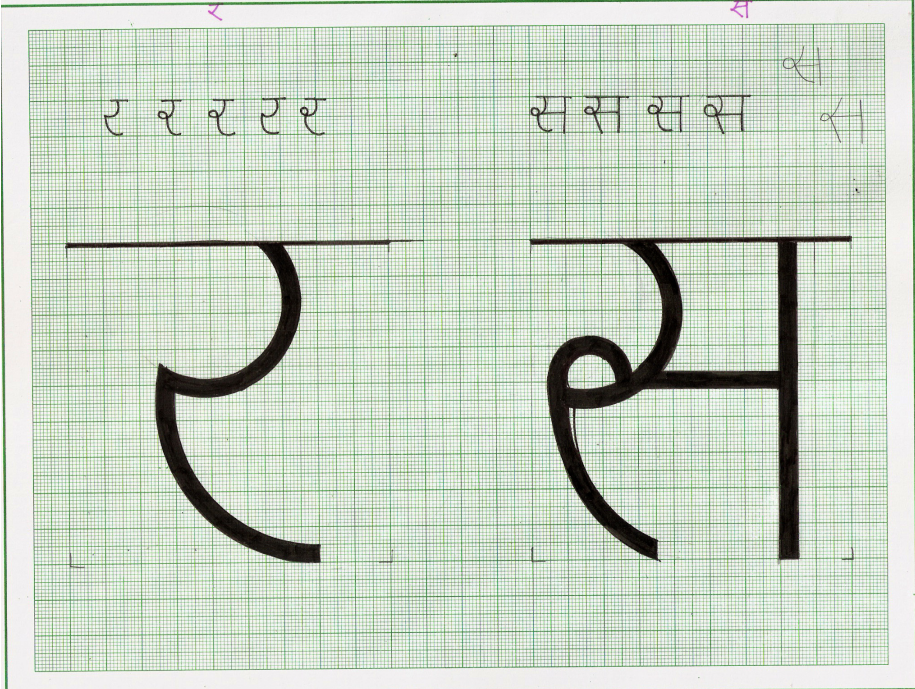


Character design; graph paper

Characters based on the grid were made on graph paper



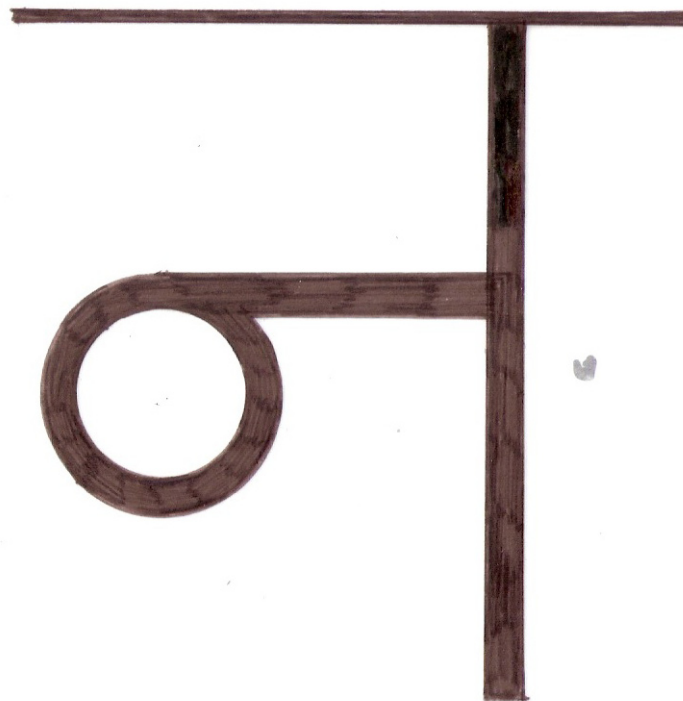
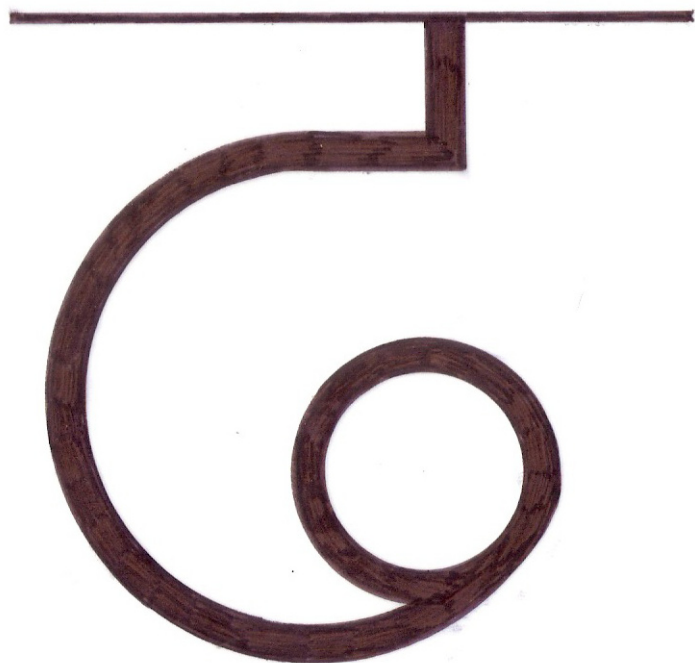
Character design; graph paper



All implicit consonants have been drawn on graph sheets

Character design; tracing sheets

All characters were traced on to sheets to be scanned for digitization



Font design and generation

- Software used for initial test typeface: Fontographer
- A new set of algorithm and codes will be written by the OCR team and the font will be tested simultaneously
- The same grid has been followed and corrections made to the character set include:
 - weight of stroke is increased
 - The distance between shirorekha and hanging characters increased
- Test prints have been reviewed and appropriate corrections have been made
- The first version of the truetype OCR font for devanagari has been sent for testing. The results of which are awaited.

Font tests

4.06.09

4.0ekimocr2.ttf. size .140pt.

इ ई ऋ क ख ग

COUNTER.

→

घ ङ ज झ ट

LOOKS CONDENSED
INCREASE WIDTH

INCREASE
LOOKS CONDENSED

ॠ ड व त द

INCREASE X

SLIGHT
STRAIGHT
GUIDE

FLAT FOOT
(refer test. 3
in smaller font size
looks limp)

CORRECT
CURVE

Font tests

4.06.09

4.06kimocr2.tff - 140pt.

प ब य र व ष

LOOP BIGGER *

स ह न न

TRY VARIATION
WITHOUT LOOP



COUNTER VARIATION
LIKE IN स

MORE OPEN
IN SMALLER
POINT SIZE

- MAKE CURVES FINE
- ABRUPT JOINERIES - CORRECT

Font tests

5.06kimocr.ttf

ह

ह
ह
ह
ह

EXTRA WEIGHT ON LEFT. - TRY MOVING/ENDING * STROKE TO RIGHT

ह

स

क
TRY OPENING LOOP
SLIGHT BIGGER LOOP

ख

ग

घ

ङ

ज

झ

ट

ठ

ड

ढ

ण

त

ह
ह
ह
ह

Font tests

Kim.Ocr.set.ttf

व य र व ष स

UNEVEN
STROKE
CHANGE TO ष

ह ा न प ँ

widen
WIDEN

COULD
CLASH
WITH
ऌ

X
प LOOKS BETTER

इ ि ि ि अ ः

REWORK.

Font tests

Font tests in various point sizes

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

इईएकखगघडझटंडढतढननपबयरवषसह

Font tests

अकखग घ ङ जझट ठ ड ढ

त द न प व य व ष स ह ङ इ भ

म ल ऐ ण प ष च उ ऊ

Font tests

अकखग घडजट ठड ढत दनपबय वषस हा इभेछच उऊत दनपबय व इभेमलपेणपयवषस हा इभेमल
मलपेणपछच उअकखग घ डजझट ठड ढत दनपबणपछच उऊत दनपबय वषस हा इभेपछच उअकखग घ ड
मलपेणपछच उट ठड ढत दनपबय वषस हा इभेणपछच उऊत दनपबय वषस हा इभेपय वषस हा पय वषस हा
मलपेणपछच उऊत दनपबय वषस हा इभेणपछच उऊत दनपबय वषस हा इभेपछच उअकखग घ ड
मलपेणपय वषस हा इभेमलपेणपय वषस हा इभेमलपेणपय वषस हा इभे

कमल नमन कर
अमर पठण कर
सहज हवन
टालमण भला यहरझ जटगयघ उहरझरे खपय हरज गपडहडट इयालरिझिपज

कमल नमन कर
अमर पठण कर
सहज हवन
टालमण भला यहरझ जटगयघ उहरझरे खपय हरज गपडहडट इयालरिझिपज

कमल नमन कर
अमर पठण कर
सहज हवन
टालमण भला यहरझ जटगयघ उहरझरे खपय हरज गपडहडट इयालरिझिपज

कमल नमन कर
अमर पठण कर
सहज हवन
टालमण भला यहरझ जटगयघ उहरझरे खपय हरज गपडहडट
इयालरिझिपज

कमल नमन कर
अमर पठण कर

सहज हवन

टालमण भला यहरझ जटगयघ उहरझरे खपय हरज
गपडहडट इयालरिझिपज

कमल नमन कर

अमर पठण कर

सहज हवन

टालमण भला यहरझ जटगयघ उहरझरे
खपय हरज गपडहडट इयालरिझिपज

कमल नमन कर

अमर पठण कर

सहज हवन

Future prospects

- Develop the entire character set, including conjuncts, explicit consonants, numerals and punctuations
- Review results from OCR team and revert back with corrections till desired quality is achieved
- Generate a devanagari font that gives optimum accuracy results for character recognition when scanned
- Proposal to continue project as Project 2