#### • ENCCS and Vega in Support of Large-Scale Digitization

- of Archival Materials
- $\bullet$

- •
- •
- •



# Digitization of Archival Data

- The Swedish National Archives has about 80000 shelf meters of archival material, ranging from the 11th century until the present
- Our digital archives consists of about 215 million scanned images of archival documents.
- Just scanning the documents don't make them searchable. For running-text documents HTR/OCR is needed for a full-text search, and for forms and tables, segmentation and HTR/OCR is needed to transform the scanned forms into searchable databases to be used by us internally, or by research



# A Few Notes on OCR/HTR

- It consists of segmentation/text-detection plus text-recognition
- In 2015 an architecture for text-recognition called CRNN (Convolutional Recurrent Neural Network) was introduced in a paper called "An End-to-End Trainable Neural Network for Image-based SequenceRecognition and Its Application to Scene Text Recognition." This architecture is still widely used, for example in Tesseract and ABBY for OCR, and by Transkribus for HTR.
- Since 2019 several architectures for text-recognition based on attention and the tranformer architecture, rather than recurrence, have been developed in research and implemented in open-source repositories such. The usual advantage of attention vs recurrence, that is, handling long-range context, also applies to text-recognition.
- Which architecture to choose is largely determined by the type of material you're going to analyze

**Riksarkivet** 

#### Indexing Projects at The Swedish National Archives

- The idea behind what we call indexing projects is to work with certain archival records that gets a lot of requests from the public, for example personal records or property records, and index them by some unique information, for example personal registration number or property name.
- By doing this the scanned images becomes searchable by this information, and the relevant personal record or property record can be quickly looked up.
- This streamlines our internal case handling and enhance the availability of the records



# Prototype Project Started in 2021

- The goal with this project was to investigate whether there is a general method and workflow for automatically indexing archival records with the help of machine-learning methods
- Our case-study was the church archive's personal records
- We didn't have a ML-infrastructure in place, and we didn't have access to any HPC-environment, so we trained models on relatively small datasets, and analyzed just a small sample of the records, with the goal of arriving at a proofof-concept
- The results were promising and it led us to hypothesize a general method for similar indexing projects at the National Archives

) )	1 FKE	71/2							and the second	Lun	d	1	M	(1000	Chi	
	Trke Fil.stud.						[ <u>]), yy</u> ]_(				A Fud-1	Malm. Hanis				
							· ·			3) Ron	4) Fodel	1 10	ocn ined	dag	1 "	
_										6) Vacc.	1926	j	uli	30	4	
8)	Födelsehemort (födort) (se ovan)						i (se ovan) länistad 7) Stam, nation.			9) Lyte		4124				
10)	Trosbekännelse							11) dopt area awker see #3-08			12) Konfirm. m. m. a. n					
13)	Börd m. m.						14) Reg. å sjömanshus			15) Värnpl.						
16)						-	-	efternamn, förnamn och födelsehemort				föselsetid och numm			umm	
drar	Fader											år	mån.	dog		
Föräl	Moder													-	1	
17)	nr		ingånge	1	. ups	trök	makes efternamn, förnamn och födelsehemort					fed	elsetid	och n	vmm	
d o		Ar K	mán.	dag	år im	àn.   d	dag'	Svensson, Axel Gunnar	and the second	84- 1		år	mån.	dag	T	
	1	1 49 07 01					Lund, Malm			:	14	08	11	41		
n s k																
Akte	2															
	-			-		+	-	14 C		1.4					t	
	3					-	-		and the second s				1			
18)	kön födelsetid och nummer		Τ.	T. A.	A. namn (övre raden) och födelsehemort (undre raden) den andras av förf		ildrarna namn		födelsetid och nur		umm					
		år mån. dag nr				Sum Transa Christenhen	So mm 17	×	×	år i	mån.	dag	1 *			
	M 51 jan 197				Sven ingemar Unristopher Se rum 17											
	-				07.1	+		Turislorp Main				1		1	T	
	K	53	dec	22	9346			Ingrid Bodil Elisabet	Fum 17	1						
		-	-	-		+	-	Tullstorp			-	+		-	t	
Barn							-									
	-			-		-	-				-	-		-	1	
8							-			· · · · · · · · · · · · · · · · · · ·						
8 a l				_		-						_		-	1	
Bar						-								1	1	
Ba	-						-		-					1		
Ba	_						_		1							
Ba	*		-					· · · · · · · · · · · · · · · · · · ·					÷			

Example of a personal record (The project manager Catharina Dahlgren's grand-mother, no GDPR problem)



#### General indexing pipeline and workflow





#### Implementing the pipeline in a production project

- Many requests from the the public made to the Swedish National Archives concerns property records. These requests requires a lot of manual work by archivists, which could be greatly reduced if the scanned images where indexed by property name
- We have approximately 18 million scanned images of property records
- About 9 million of them are already manually indexed during a period of ten years, which means we have all the training data we need.



#### The Property Record Indexing Pipeline





#### Label Studio and MMDetection for Object Detection

Label Studio = Projects / property_detection_40_batches_train_700_crop / Labeling	Label Studio       E       Projects / property_detection_40_batches_train_700_crop / Labeling
Kvärteret Almen tomten nr 6       Råsunda mie Stadsdel     Stadsdel	Hosjöholmen <sup>Blad #</sup> 6'
Solna Stad Fastighetsespalt. 109%	property 1 Eabel Studio E Projects / property_detection_40_batches_train_700_crop /
property 1 $\begin{array}{c c} Label Studio \\ \hline \\ $	Projects / property_detection / Labeling
III Appart <ul> <li>Inlogging - starm</li> <li>I Label Studio</li> <li>Projects / property_detection_40_batches_train_700_crop / Labeling</li> </ul>	property 1
Label Studio Projects / property_detection_40_batches_train_700_crop / Labeling	Projects / property_detection_40_batches_train_700_crop / Labeling
Jomlen mr. J. ko. Vaktaren. Intecknings-spalt. 218	Label Studio = Projects / property_detection_40_batches_train_700_crop / Labeling
property 1	Kwarteret tomten nr 5 Almen Stadsägan nr 16.76.
Riksarkivet	Spalt för antockning om fastighetens natur 2232 Sätunder som hans fännen. Den då lanfort

#### Training the Text-Recognition Model

LANKOD; KOMKOD; FNR; GKOMMUN; GTRAKT; GBLOCK; GENHET; DATUM; AKT; KOMMUN; TRAKT; RBLOCK; ENHET 01;88;010000002;A-GOTTRÖRA;ABRAHAMSBY;1;1;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;1;1 01;88;010000003;A-GOTTRÖRA;ABRAHAMSBY;2;1;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;1 01;88;010000004;A-GOTTRÖRA;ABRAHAMSBY;2;2;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;2 01;88;010000005;A-GOTTRÖRA;ABRAHAMSBY;2;3;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;3 01;88;010000006;A-GOTTRÖRA;ABRAHAMSBY;2;4;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;4 01;88;010000007;A-GOTTRÖRA;ABRAHAMSBY;2;5;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;5 01;88;010000008;A-GOTTRÖRA;ABRAHAMSBY;2;6;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;6 01;88;010000009;A-GOTTRÖRA;ABRAHAMSBY;2;7;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;7 01;88;010000010;A-GOTTRÖRA;ABRAHAMSBY;2;8;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;8 01;88;010000011;A-GOTTRÖRA;ABRAHAMSBY;2;9;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;9 01;88;010000012;A-GOTTRÖRA;ABRAHAMSBY;2;10;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;10 01;88;010000013;A-GOTTRÖRA;ABRAHAMSBY;2;11;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;11 01;88;010000014;A-GOTTRÖRA;ABRAHAMSBY;2;12;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;12 01;88;010000015;A-GOTTRÖRA;ABRAHAMSBY;2;13;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;13 01;88;010000016;A-GOTTRÖRA;ABRAHAMSBY;2;14;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;14 01;88;010000017;A-GOTTRÖRA;ABRAHAMSBY;2;15;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;15 01;88;010000018;A-GOTTRÖRA;ABRAHAMSBY;2;16;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;16 01;88;010000019;A-GOTTRÖRA;ABRAHAMSBY;2;17;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;17 01;88;010000020;A-GOTTRÖRA;ABRAHAMSBY;2;18;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;18 01:88:010000021:A-GOTTRÖRA:ABRAHAMSBY:2:19:19791001:01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2:19 MMOCR, SATRNlamten mi 3 a praiteret bland OJABY 14 Architecture Öjaby Persgård 6:261. Öjaby Getaskärv Norregård 242 Gemla 5:88 Tomten nr 3 i kvarteret Romaren. Ellanda 7<sup>+</sup>• Using our object-detection model, and the database for Øjaby Ingelsgård 9:117 the already indexed property records we created the training set for the text-recognition model. We trained it on Vega, on approximately one million images Getaskarv Norregard 2:42 Östanstorp 31. sarkivet

HTR/

OCR-

model

#### Why Attention Rather than Recurrence?



![](_page_10_Picture_2.jpeg)

#### **Post-Correction and Data Validation**

LANKOD; KOMKOD; FNR; GKOMMUN; GTRAKT; GBLOCK; GENHET; DATUM; AKT; KOMMUN; TRAKT; RBLOCK; ENHET 01;88;010000002;A-GOTTRÖRA;ABRAHAMSBY;1;1;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;1;1 01;88;010000003;A-GOTTRÖRA;ABRAHAMSBY;2;1;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;1 01;88;010000004;A-GOTTRÖRA;ABRAHAMSBY;2;2;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;2 01;88;010000005;A-GOTTRÖRA;ABRAHAMSBY;2;3;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;3 01;88;010000006;A-GOTTRÖRA;ABRAHAMSBY;2;4;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;4 01;88;010000007;A-GOTTRÖRA;ABRAHAMSBY;2;5;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;5 01;88;010000008;A-GOTTRÖRA;ABRAHAMSBY;2;6;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;6 01;88;010000009;A-GOTTRÖRA;ABRAHAMSBY;2;7;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;7 01;88;010000010;A-GOTTRÖRA;ABRAHAMSBY;2;8;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;8 01;88;010000011;A-GOTTRÖRA;ABRAHAMSBY;2;9;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;9 01;88;010000012;A-GOTTRÖRA;ABRAHAMSBY;2;10;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;10 01;88;010000013;A-GOTTRÖRA;ABRAHAMSBY;2;11;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;11 01;88;010000014;A-GOTTRÖRA;ABRAHAMSBY;2;12;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;12 01;88;010000015;A-GOTTRÖRA;ABRAHAMSBY;2;13;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;13 01;88;010000016;A-GOTTRÖRA;ABRAHAMSBY;2;14;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;14 01;88;010000017;A-GOTTRÖRA;ABRAHAMSBY;2;15;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;15 01;88;010000018;A-GOTTRÖRA;ABRAHAMSBY;2;16;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;16 01;88;0100000019;A-GOTTRÖRA;ABRAHAMSBY;2;17;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;17 01;88;010000020;A-GOTTRÖRA;ABRAHAMSBY;2;18;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;18 01;88;010000021;A-GOTTRÖRA;ABRAHAMSBY;2;19;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;19 01;88;010000022;A-GOTTRÖRA;ABRAHAMSBY;2;20;19791001;01-BJÖ-836;NORRTÄLJE;ABRAHAMSBY;2;20

We validated the predictions against an existing database of property names, to make sure the prediction is a valid property name. No guarantee against false positives. Various other validation criteria were also applied

![](_page_11_Picture_4.jpeg)

If validation fails, the image gets flagged for manual control in an indexing program developed for this specific purpose

![](_page_11_Picture_6.jpeg)

# Manual Correction for Non-Validated Predictions

Manual correction where:

- Link agains validation database fails
- Low prediction-confidence
- Valid property name, but wrong region

Tradeoff between number of false positives and number of hits.

00000405 203 2 Fill		Dubblett		Klar för kontroll 🛛 🗸	10005399 ~
ÄLMHULT	@:RENEN:8			back next Ko	Välj alla batcher 1 zo
<u>Xlmhul t</u>	Renen 8	Arkivblad B	<i>T 15900</i> lad 1	2 1	Start Sök
	Fastighets-s	palt		1 @;RENEN;8 2 @;RENEN;8	MARKARYD Älmhult
	Renen 8 Lagfarts-spa	1 t		check_Refere	enskod_serie_id Nästa Batch
00000405 4 203 2 1256256 070110343	ALMHULT @;RENEN;8	- Dubblett	K.före K.efter ??? Inf	fo 3	Facit
00000407 2 204 0 1256256 070110343	ALMHULT @;RENEN;8	- Dubblett	K.före K.efter ??? Inf	fo 50 3	0 Facit
00000409 2 205 0 1256259 070110344	ÄLMHULT @/RENEN:9	- Dubblett	K.fore K.efter ??? Inf		0 Facit
	LUNGBY @:FÖRTENNAREN:5	- Dubberr	K.före K.efter ??? Inf		Large + (OK) Eacit
00000415 2 208 1 1 1256636 070110709	STENBROHULT KVARNATORP:1:35		K före K efter 222 Inf		Large + Facit
00000417 2 209 2 1 1238894 070095823	MARKARYD @:STG:347	-	K.före K.efter ??? Inf	fo 1 1	Large + (OK) Facit
00000419 0 210		-	K.före K.efter ??? Inf	fo 0 0	Missing Facit
Next         180         Dubblett         Sant         2           216         Dubblett         Falskt         2           230         Dubblett         Falskt         2	1 203 @;RENEN;8 2 203 @;RENEN;8	MARKARYD 1 0 ÄLMHULT 1 0	070097643 070110343	2 Clear 2 Skapa gemensa	umt index. Visa gemensamma inde
Previous         222         Dubblett         Falskt 2           21         Dubblett         Falskt 2         Falskt 2           232         Dubblett         Falskt 2	< <			> omg1 om	g2 omg3 Beräkna om stat

![](_page_12_Picture_7.jpeg)

#### Results

- We evaluated the pipeline on a testset of 100000 images spread over the data-space (locations in Sweden where the records were made). We got exact matches on just over 90% of the property names.
- Below is a venn diagram of the 233 batches the testset consisted of. Each batch consists of roughly 500 images. The bars represent 5 percent intervals. Most of the batches below 70% were actually errors in ground truth, that is, errors in the manual indexing.

![](_page_13_Figure_3.jpeg)

![](_page_13_Picture_4.jpeg)

## Vega for training and inference

- Training the SATRN-model on Vega enabled us to increase the scale of the resized images going into the model, thereby improving accuracy for handwritten text, which generally requires more information than printed text
- Running 9 million images thorugh the pipeline on VEGA took roughly 90 node-hours
- At a hit-rate of 90% this project saves us about 700000 euros in manual labor costs, and the indexing database gets created a lot quicker

![](_page_14_Picture_4.jpeg)

#### ENCCS and Vega are opening up huge possibilities

- More of the same, indexing projects that saves tax-money and streamlines our casehandling
- We are also in the process of creating running-text HTR base-models for different timeperiods with the goal of mass-digitizing historical handwritten archival documents and making them available for research and for the public
- We want to improve our OCR-pipeline and apply it on scale
- We would also like to create base-models for segmentation of historical forms and tables
- One possible future project is to automatically transcribe the 19th-century censuses, this project would be a part of an already existing research infrastructure called swed-pop, of which the National Archives is a part, and where we are doing this transcriptions manually
- Inspired by the Royal Library of Sweden, we've also done experiments with training historical LM:s, that is, language models like BERT, adapted to historical text.

![](_page_15_Picture_7.jpeg)

# : Thank you for listening,

- And big thanks to:
- ENCCS,
  - EuroHPC,
  - and Vega

![](_page_16_Picture_5.jpeg)