

Constructing a Historical Nordic Human Capital Database: An End-to-End Machine Learning Approach

Christian E. Westermann & Christian M. Dahl

University of Southern Denmark

December 31, 2021

Motivation

- A unique database on human capital in the Nordic countries, spanning multiple centuries.
- With successful links to current registries, large potential for multi-generational analyses.
- To the best of our knowledge, no end-to-end demonstration of historical tabular data segmentation and transcription exists.
- Structured tabular data are extremely valuable, especially for economists (census data, etc.).
- Many subtle but extremely important complications.
- Contribute a generalizable approach with respect to structured document tables.

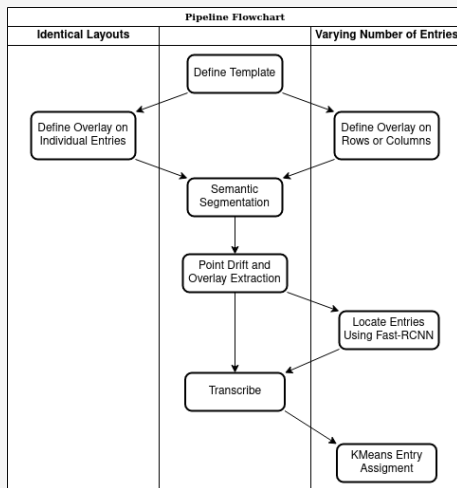
Related Work

- Many parallels to dhSegment, the work of Oliveira et al. (2019).
 - They focus on deep learning based generic segmentation of several tasks such as page-, image-, and text-detection.
 - We wish to introduce a generic approach to structured tabular data.
- LayoutParser by Shen et al. (2021).
 - Unified toolkit for Deep Learning based DIA.
 - Includes layout detection, OCR with Google Vision or Tesseract backend and more.
- Google's Tesseract
 - While sufficient in many tasks, it is too inaccurate for the tables we work with.
 - Too many or not enough boxes detected.
 - Very hard to sort relevant boxes.

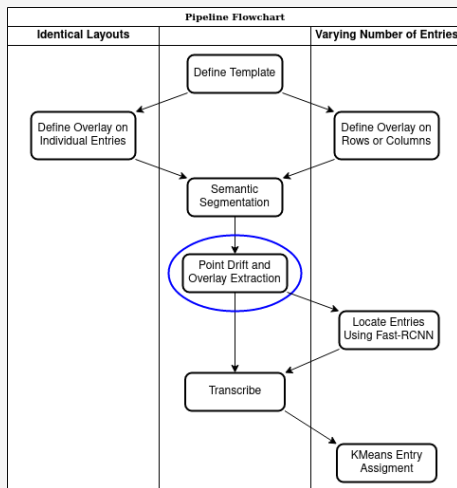
Tesseract Example

Fortegnelse ordnet				Efter Eksamensstedet											
Afskriftet af				Afskriftet af											
Da Staderendes Navne.				Afskriftet af											
Metropolitanskolen.				Afskriftet af											
Richter, Vilhelm				Afskriftet af											
Bache, Niels				Afskriftet af											
Kobke, Hans Peter Carl				Afskriftet af											
Sahlert, Ivan Edvard (Alexander)				Afskriftet af											
Pach, Theodor Emil Bartholomæus				Afskriftet af											
Krohn, Peter Jakob				Afskriftet af											
Hjorth, Hans Christian Carl Mønst.				Afskriftet af											
Larsen, Vilhelm Christian Sigurd				Afskriftet af											
Hansen, Peter Guillelmus				Afskriftet af											
Nielsen, Albert Emil Adolph				Afskriftet af											
Wiise, Georg Carl Christian				Afskriftet af											
Bender, Viggo Henrik Lauritz				Afskriftet af											
Kier, Edgar Vilhelm				Afskriftet af											
Heerfordt, Niels Christian				Afskriftet af											
Thølle, Henrik Theodor Jacob				Afskriftet af											
Roeskilde Skole.				Afskriftet af											
Gentzen, Henrik Michael				Afskriftet af											
van Wylich, William Gustav				Afskriftet af											
Breyer, Jens Nicolai				Afskriftet af											
Bagger, Otto Mandix				Afskriftet af											
Holstein-Lethborg, Johan Ludvig Carl Christian				Afskriftet af											
Lund, Laurits Christian				Afskriftet af											
Lund, Emmannuel Johannes Christian				Afskriftet af											
Frederiksberg Skole.				Afskriftet af											
Møller, Christian				Afskriftet af											
Blichfeld, Michael Frederik Kamman				Afskriftet af											
Horsp, Goshald Carl (Anders) Oluf				Afskriftet af											
Udall, Carl Julius				Afskriftet af											
Holmblad, Jacob Arnold Christian				Afskriftet af											

Approach Overview



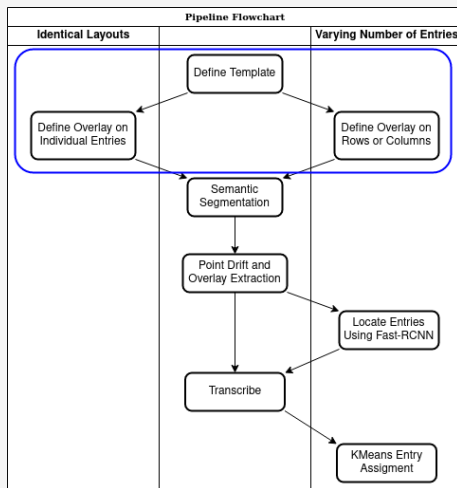
Approach Overview



Point Drift

- The engine behind the approach is FilterReg by Wei Gao and Russ Tedrake (2018), a Point-Set Registration algorithm.
- Learn the Motion (Transformation) Parameters, $\Delta\theta$, responsible for the alignment of two point clouds, X and Y .
- We can then apply $\Delta\theta^{-1}$ to Y , which aligns it with X .

Approach Overview



Template and Overlay

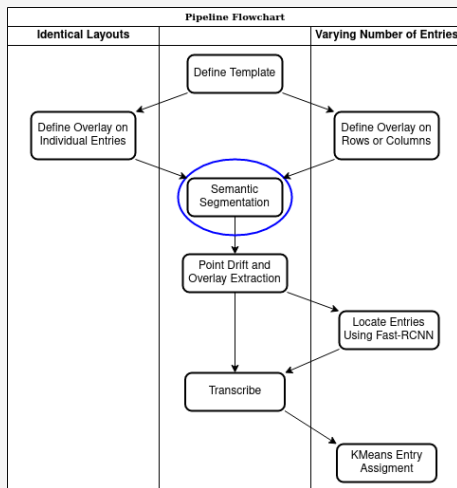
- Manually define two setup steps per document layout type: template and overlay.
- We define the template as the document table outline.
- How to define the overlay is dependent on whether number of entries vary.
 - Static: Entry-level overlay
 - Varying: Column- or Row-level overlay.

Template and Overlay

Nævne	Skole eller Bistands.	Præ- dødt	legit.	Psyke- legit.	Fødsels- dag
Jensen, L. B.	Fredriksholm Skole.	gødt.	mg.	gødt.	23. Juni.
Jørgensen, J. V.	Randers Skole.	mg.	mg.	mg.	15. Juni.
Johnsen, J. V.	Metropolitanskolen.	mg.	mg.	mg.	18. Juni.
Johnsen, O. H.	Reykjavik Skole.	mg.	mg.	mg.	18. Juni.
Jørn, H. A.	Borgerskole, paa Chavn.	mg.	mg.	mg.	15. Juni.
Jørgensen, C.	Flebens Skole.	mg.	mg.	mg.	16. Juni.
Kiats, P. J. W.	Herfsholms Skole.	gødt.	gødt.	gødt.	24. Juni.
Kleisborg, G. M.	Ribe Skole.	mdl.	mg.	mg.	18. Juni.
Koch, H. L. S. P.	Nykjøbing Skole.	mg.	gødt.	mg.	10. Juni.
Krup, C. F. E.	Flebens Skole.	gødt.	mg.	mg.	12. Juni.
Krup, H. A.	Borgerskole, paa Chavn.	tg.	mg.	mg.	23. Juni.
Landsgjort, Th. Ph.	Fredriksholm Skole.	mg.	gødt.	mg.	19. Juni.
Langkilde, F. E.	Odense Skole.	mg.	mg.	mg.	23. Juni.
Larsen, A. C.	Aalborg Skole.	gødt.	mg.	mg.	9. Juni.
Larsen, H. C. A.	Fredriksholm Skole.	mg.	mg.	mg.	20. Juni.
Laury, C. L. B.	Aarhus Skole.	tg.	mg.	mg.	31. Januar.
Leth, C. P.	Soro Skole.	mg.	mg.	mg.	16. Juni.
Leth, J. Q.	Aalborg Skole.	gødt.	mg.	mdl.	31. Januar.
Leunbach, H. G.	Horsens Skole.	mg.	mg.	mg.	10. Juni.
Lorenzen, C. N.	Flebens Skole.	gødt.	mg.	mg.	20. Juni.
Lochte, H. N. J.	Aalborg Skole.	gødt.	mg.	mg.	30. Januar.
Madsig, P. A. G.	Conferens. Mødt.	gødt.	tg.	13. Juni.	
Martensen, C.	Metropolitanskolen.	mg.	mg.	mg.	20. Juni.
Martensen, C. J.	Sanne Skole.	mg.	mg.	mg.	20. Juni.
Mohr, J. J.	Metropolitanskolen.	mg.	mg.	mg.	24. Juni.
Mohr, S. J. G.	Aarhus Skole.	mg.	gødt.	gødt.	27. Juni.
Müller, P. G.	Ribe Skole.	mg.	gødt.	gødt.	17. Juni.
Müller, J. J. N.	Roskilde Skole.	mg.	mg.	mg.	16. Juni.
Nissen, N. C. A.	Cand. ph. C. Kofod.	gødt.	gødt.	gødt.	20. Juni.
Norregaard, J.	Borgerskole, paa Chavn.	mg.	mg.	mg.	10. Juni.
Olivarius, H. H. F.	v. Westenske Institut.	tg.	gødt.	mg.	23. Juni.
Olrik, H. L. Th.	Herfsholms Skole.	mg.	mg.	mg.	19. Juni.
Paulsen, L. C. C.	Cand. jmr. Nolleman.	gødt.	mg.	mg.	30. Januar.
Petersen, C.	Borgerskole, paa Chavn.	gødt.	gødt.	tg.	27. Juni.
Petersen, A. N. Fallman.	Soro Skole.	gødt.	mg.	gødt.	22. Juni.
Petersen, C. H. W.	Haderslev Skole.	gødt.	gødt.	mg.	15. Juni.
Petersen, P. A.	Cand. ph. B. Møller.	mg.	mg.	mg.	9. Juni.
Petersen, R.	Haderslev Skole.	gødt.	gødt.	mg.	30. Januar.

Nævne	Skole eller Bistands.	Præ- dødt	legit.	Psyke- legit.	Fødsels- dag
Jensen, L. B.	Fredriksholm Skole.	gødt.	mg.	gødt.	23. Juni.
Jørgensen, J. V.	Randers Skole.	mg.	mg.	mg.	15. Juni.
Johnsen, J. V.	Metropolitanskolen.	mg.	mg.	mg.	18. Juni.
Johnsen, O. H.	Reykjavik Skole.	mg.	mg.	mg.	18. Juni.
Jørn, H. A.	Borgerskole, paa Chavn.	mg.	mg.	mg.	15. Juni.
Jørgensen, C.	Flebens Skole.	mg.	mg.	mg.	16. Juni.
Kiats, P. J. W.	Herfsholms Skole.	gødt.	gødt.	gødt.	24. Juni.
Kleisborg, G. M.	Ribe Skole.	mdl.	mg.	mg.	18. Juni.
Koch, H. L. S. P.	Nykjøbing Skole.	mg.	gødt.	mg.	10. Juni.
Krup, C. F. E.	Flebens Skole.	gødt.	mg.	mg.	12. Juni.
Krup, H. A.	Borgerskole, paa Chavn.	tg.	mg.	mg.	23. Juni.
Landsgjort, Th. Ph.	Fredriksholm Skole.	mg.	gødt.	mg.	19. Juni.
Langkilde, F. E.	Odense Skole.	mg.	mg.	mg.	23. Juni.
Larsen, A. C.	Aalborg Skole.	gødt.	mg.	mg.	9. Juni.
Larsen, H. C. A.	Fredriksholm Skole.	mg.	mg.	mg.	20. Juni.
Laury, C. L. B.	Aarhus Skole.	tg.	mg.	mg.	31. Januar.
Leth, C. P.	Soro Skole.	mg.	mg.	mg.	16. Juni.
Leth, J. Q.	Aalborg Skole.	gødt.	mg.	mdl.	31. Januar.
Leunbach, H. G.	Horsens Skole.	mg.	mg.	mg.	10. Juni.
Lorenzen, C. N.	Flebens Skole.	gødt.	mg.	mg.	20. Juni.
Lochte, H. N. J.	Aalborg Skole.	gødt.	mg.	mg.	30. Januar.
Madsig, P. A. G.	Conferens. Mødt.	gødt.	tg.	13. Juni.	
Martensen, C.	Metropolitanskolen.	mg.	mg.	mg.	20. Juni.
Martensen, C. J.	Sanne Skole.	mg.	mg.	mg.	20. Juni.
Mohr, J. J.	Metropolitanskolen.	mg.	mg.	mg.	24. Juni.
Mohr, S. J. G.	Aarhus Skole.	mg.	gødt.	gødt.	27. Juni.
Müller, P. G.	Ribe Skole.	mg.	gødt.	gødt.	17. Juni.
Müller, J. J. N.	Roskilde Skole.	mg.	mg.	mg.	16. Juni.
Nissen, N. C. A.	Cand. ph. C. Kofod.	gødt.	gødt.	gødt.	20. Juni.
Norregaard, J.	Borgerskole, paa Chavn.	mg.	mg.	mg.	10. Juni.
Olivarius, H. H. F.	v. Westenske Institut.	tg.	gødt.	mg.	23. Juni.
Olrik, H. L. Th.	Herfsholms Skole.	mg.	mg.	mg.	19. Juni.
Paulsen, L. C. C.	Cand. jmr. Nolleman.	gødt.	mg.	mg.	30. Januar.
Petersen, C.	Borgerskole, paa Chavn.	gødt.	gødt.	tg.	27. Juni.
Petersen, A. N. Fallman.	Soro Skole.	gødt.	mg.	gødt.	22. Juni.
Petersen, C. H. W.	Haderslev Skole.	gødt.	gødt.	mg.	15. Juni.
Petersen, P. A.	Cand. ph. B. Møller.	mg.	mg.	mg.	9. Juni.
Petersen, R.	Haderslev Skole.	gødt.	gødt.	mg.	30. Januar.

Approach Overview



Semantic Segmentation

- Now that we have **X**, how do we find **Y**?
- Semantic segmentation (pixel-wise classification) using Deep Learning.
- Three classes: Line, text / scribbles and background.
- Pixels that make up detected lines become the **Y** point cloud.

Semantic Segmentation: Neural Network Architecture

- We use the architecture of DeeplabV3+, by Chen et al. (2018).
- U-shaped and like other U-shaped architectures since originally introduced by Ronneberger et al. in 2015, it has proven to excel at semantic segmentation.
- We train two separate models, one for horizontal line identification and one for vertical.

Semantic Segmentation: Training Data

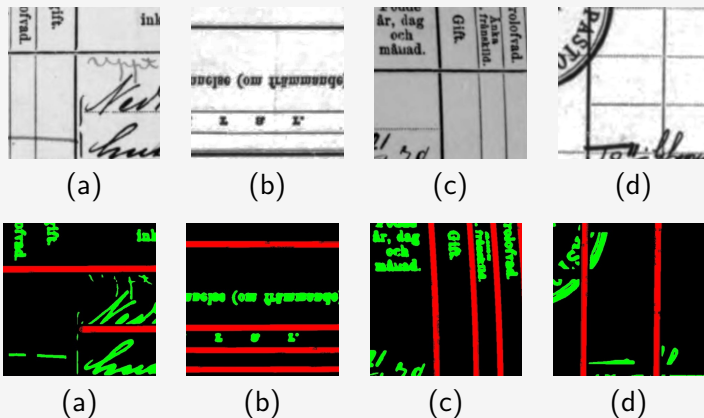
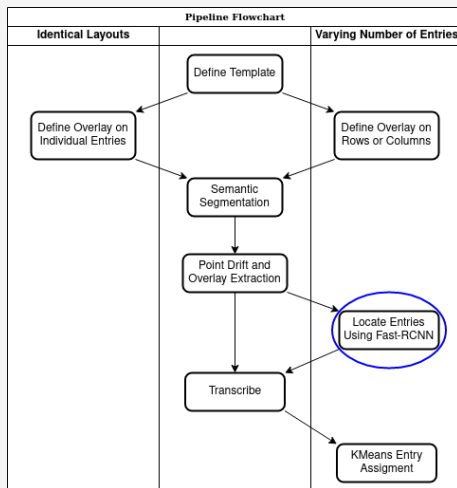


Figure: Top row: Input training images. Bottom row: Corresponding labels / masks. Red, green and black correspond to lines, text and background respectively.

15 / 22

Approach Overview



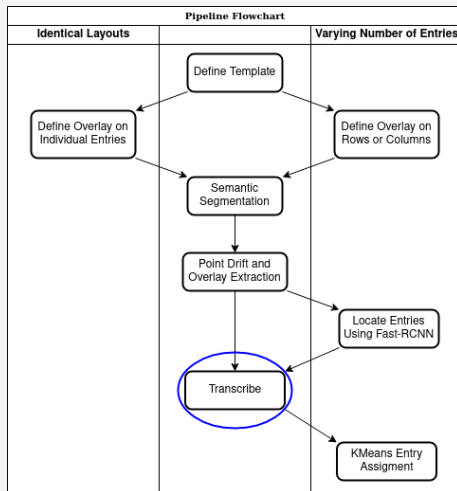
Entry Detection Using Fast-RCNN

- Faster Region-based Convolutional Network (Faster-RCNN) by Ren et al. (2016).
- Performed column-wise on extracted areas using the overlay.
- We know what column we are dealing with and we can then proceed on locating entries top to bottom.
- Network was trained on data in this paper, but we seek generality here as well, in the future.

Entry Detection Using Fast-RCNN: Results

<p>5) Fra Borgerdydskolen paa Christiansh.</p> <p>25. <i>Corfixen, Hans</i></p> <p>26. <i>Benedictsen, Boje</i></p> <p>27. <i>Bruun, Peter Adolph Rostgaard</i></p> <p>28. <i>Lorenzen, Henning Nis</i></p> <p>29. <i>Krarup-Vilstrup, Andreas Fabricius</i></p> <p>30. <i>Lund, Frederik Christian</i></p> <p>31. <i>Dresing, Frederik Nicolai</i></p> <p>32. <i>Wroblewsky, Johannes Julius</i></p> <p>33. <i>Brusch, Jens Ludvig</i></p> <p>34. <i>Flagstad, Paul Wilhelm</i></p> <p>35. <i>Buch, Edvard Magnus</i></p> <p>36. <i>Ovesen, Wolf Frederik</i></p> <p>37. <i>Schow, Ulrik Frederik Rosing</i></p> <p>38. <i>Dorscheus, Andreas Peter</i></p> <p>39. <i>Dall, Frederik Julius</i></p>	<p>Skole eller Dimissor, Ch</p> <p>Cand. th. M. T. Becher</p> <p>v. Westenske Institut</p> <p>Samme Skole</p> <p>Borgd. Skol. paa Chvn.</p> <p>Borgd. Skol. i Rbhvn.</p> <p>St. med. & chir. Trier</p> <p>Horsens Skole</p> <p>Tentamen i Roeskilde</p> <p>Bessestad Skole</p> <p>Borgd. Skol. paa Chvn.</p> <p>Samme Skole</p> <p>Samme Skole</p> <p>Stud. philol. N.E. Riis</p> <p>Aarhuus Skole</p> <p>Ribe Skole</p> <p>Frederiksborg Skole</p> <p>Borgd. Skol. i Rbhvn.</p> <p>Roeskilde Skole</p> <p>Borgd. Skol. paa Chvn.</p> <p>Samme Skole</p> <p>Cand. juris Schack</p> <p>Katechet Rasmussen</p> <p>Horsens Skole</p> <p>Roeskilde Skole</p> <p>Borgd. Skol. paa Chvn.</p> <p>Samme Skole</p> <p>Borgd. Skol. i Rbhvn.</p> <p>Roeskilde Skole</p>	<p>Arithmetik.</p> <p>Laud.</p> <p>Laud.</p> <p>Laud.</p> <p>Laud. p.c. L</p> <p>Laud. p.c. L</p> <p>H. ill.</p> <p>Laud.</p> <p>Laud. p.c. L</p> <p>Laud.</p> <p>Laud.</p> <p>Laud. p.c. L</p> <p>Laud.</p> <p>Laud.</p> <p>Laud.</p> <p>Laud.</p> <p>Laud.</p> <p>Laud.</p> <p>Laud.</p> <p>Non cont. N</p> <p>Laud.</p> <p>H. ill.</p>
---	---	---

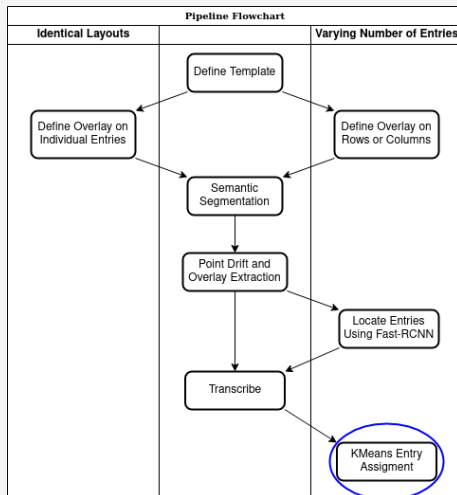
Approach Overview



Transcription

- Transcription can be done using any feasible model.
- Given the precise cutouts together with a little cleaning, Tesseract performs very well on these documents.
- Otherwise, you would have to do some fine-tuning / custom training.

Approach Overview



Assigning Entries Using KMeans

- Crucial step of the pipeline for documents with no constant number of entries.
- Missing an entry propagates the error which in the worst case leads to a full document of wrong assignments.
- We alleviate this by assigning entries using KMeans-clustering.
 - Issue a majority vote on the number of entries.
 - Create a cluster for each row based on the center y-coordinate of each bounding box.
 - Iterate through entries in each column and predict their respective row.