

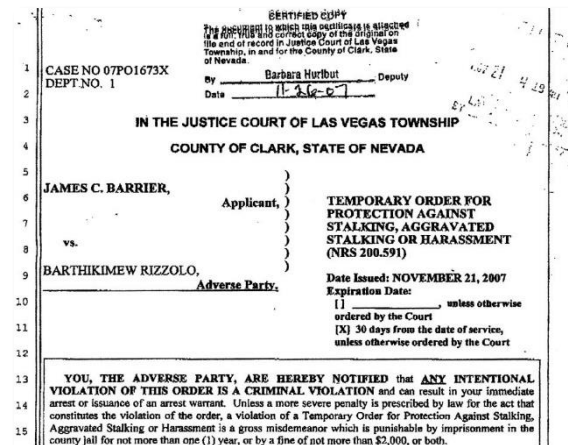
# Deep Learning für Dokumentenanalyse

*Kolloquium der ETH Bibliothek, online, 25. März 2021*

Thilo Stadelmann



# Document analysis?



## Documents

- **Ubiquitous** in human communication and every scenario involving an office
- Somewhat structured for human expert; **unstructured** w.r.t machines
- **Great use case for various AI techniques**, including computer vision

## Own scientific community

- IAPR's biannual Intl. Conference on Document Analysis & Recognition (ICDAR): character & symbol recognition, printed/handwritten text recognition, graphics analysis & recognition, document analysis & understanding, historical documents & digital libraries, document-based forensics, camera & video-based scene text analysis

# Content-Based Video Retrieval in Historical Collections of the German Broadcasting Archive

Markus Mühling<sup>1</sup>, Manja Meister<sup>4</sup>, Nikolaus Korfhage<sup>1</sup>, Jörg Wehling<sup>4</sup>, Angelika Hörth<sup>4</sup>, Ralph Ewerth<sup>2,3</sup>, and Bernd Freisleben<sup>1</sup>

The screenshot shows a web-based search interface for video content. It features a grid of video thumbnails organized into three columns: 'Anfang' (beginning), 'Mitte' (middle), and 'Ende' (end). Each thumbnail is accompanied by a timestamp. To the right of the grid is a search control panel with sections for 'Konzepterkennung' (concept recognition), 'Personenerkennung' (person recognition), 'Ähnlichkeitssuche' (similarity search), and 'OCR-Suche' (OCR search). A 'Textsuche' button is also visible at the bottom right of the interface.

# Examples

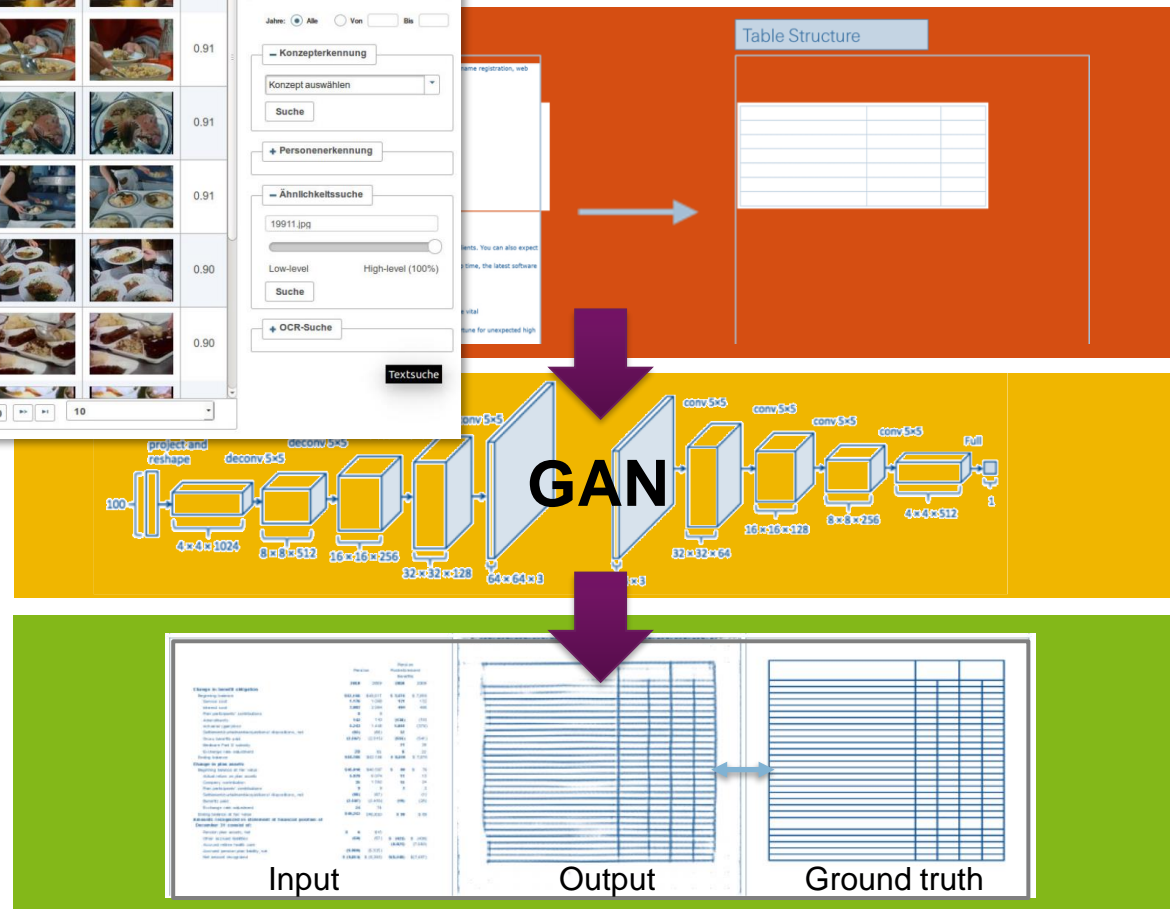
The screenshot displays the Turicode website with the headline 'Data Extraction Using Machine Learning Turn any document into structured data in just a few seconds.' Below the headline, it states: 'Stop processing your data manually – take your business process automation to a new level with the most comprehensive document understanding platform.' There are buttons for 'TURICODE SOFTWARE' and 'WATCH ON-DEMAND DEMO'. A sidebar on the right shows a list of document headers.

**Any Document Type, Any Language**  
Use our intuitive interface to process any file type; in any language, with any structure and layout.

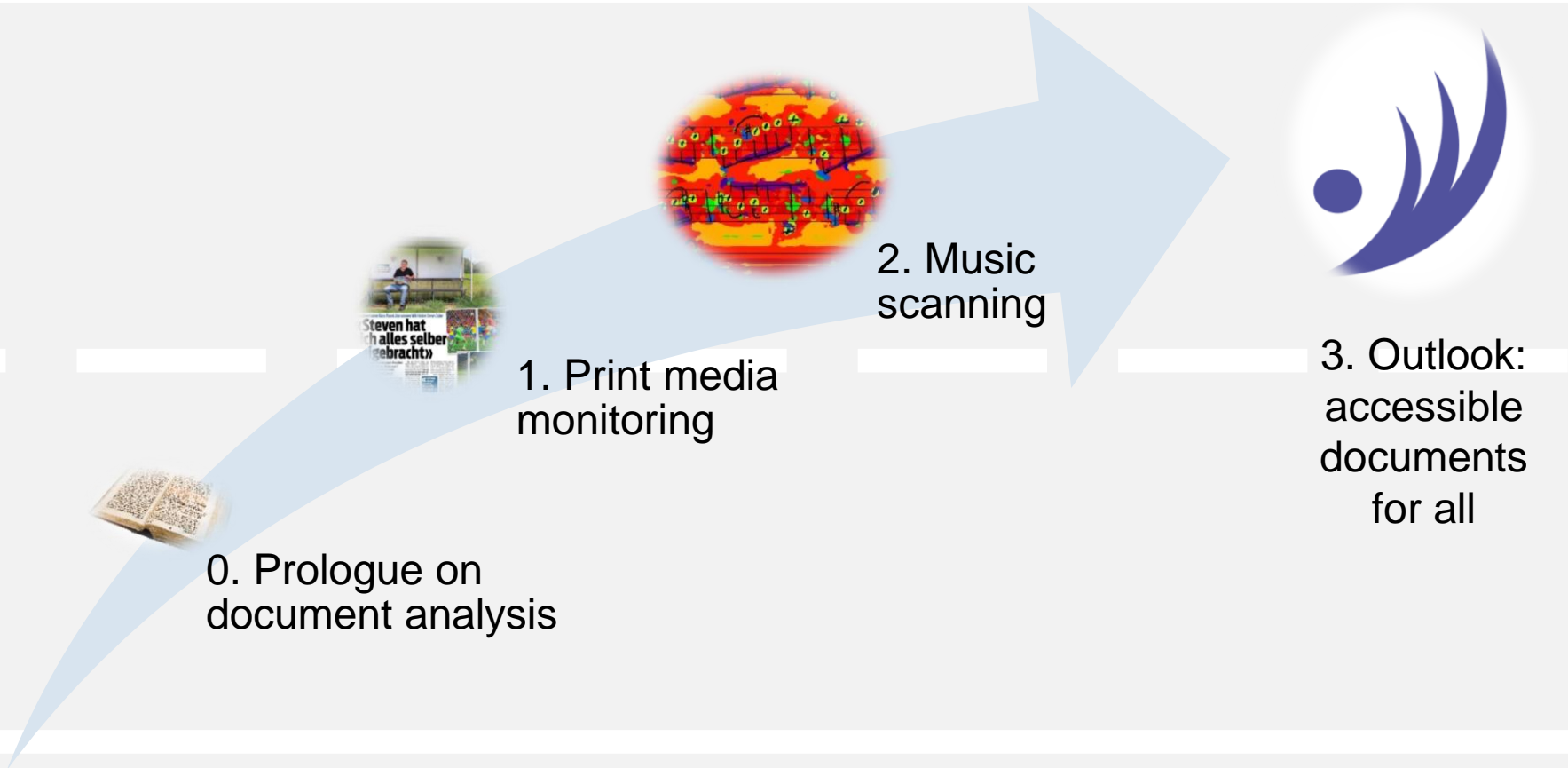
**Now Anyone Can Use Machine Learning**  
It's never been easier to train and control the AI-powered engine with just a few clicks. No coding skills required.

**Seamless Integration w Your Systems**  
Flexible data export via API to feed and enhance your ERP, CRM, or RPA system of choice.

35 Million Pages Analyzed      5 Years of Experience      40+ Document Types Automatically Analyzed



# Roadmap



# 1. Print media monitoring

## Task

International. Blick 15 | 16.05.2018, 10.44.27.18 | 1/3884

### Nachrichten

#### Spionage für den Erzfeind Iran

Iranischer Ex-Minister arbeitete als Agent für die Mullahs. Jetzt droht ihm lebenslanglich



**Amir Amirhossein** hat sich als Spion für den Iran betätigt. Der Ex-Minister arbeitete als Agent für die Mullahs. Jetzt droht ihm lebenslanglich.

Amir Amirhossein ist ein iranischer Ex-Minister, der als Spion für den Iran arbeitete. Er wurde für seine Spionagetätigkeiten verurteilt und droht ihm lebenslangliche Haft.

#### Vorbereitung

Vorbereitung der Spionagetätigkeiten.

#### Verdacht

Verdacht auf Spionagetätigkeiten.

#### Vermögen beschlagnahmt

Vermögen beschlagnahmt.

#### Asylbewerber können bleiben

Asylbewerber können bleiben.

#### Nordkoreanischer Diktator zu Besuch in Peking

Nordkoreanischer Diktator zu Besuch in Peking.

## Challenge

Sport | Blick 15 | 16.05.2018

### Sein Juniorkontainer Mano Paves über unseren WM-Helden Steven Zuber

# «Steven hat sich alles selber beibracht»

Hinter dem Zuber-Glück gegen Brasilien steckt auch Mano Paves. Mano war Der Routinedrillmeister vor dem FC Koblern im Rhein-Kreis-Stevens



Der Routinedrillmeister vor dem FC Koblern im Rhein-Kreis-Stevens. Steven Zuber hat sich alles selber beibracht.

#### Transfer Ticker

Liverpool will Yann Sommer

## Nuisance

Freitag, 22. April 2018

### Das Tages-Horoskop

**Liebling der Steine**  
**Lowe** 231-23R

<b>Jungfrau</b> 24.8.-23.9.	<b>Waage</b> 24.9.-23.10.	<b>Skorpion</b> 24.10.-22.11.	<b>Liebling der Steine</b> 231-23R
<b>Schütze</b> 23.11.-21.12.	<b>Steinbock</b> 22.12.-20.1.	<b>Wassermann</b> 21.1.-19.2.	<b>Fische</b> 20.2.-19.3.
<b>Widder</b> 21.3.-20.4.	<b>Stier</b> 21.4.-20.5.	<b>Zwillinge</b> 21.5.-20.6.	<b>Krebs</b> 20.6.-22.7.

### 15,1 Millionen

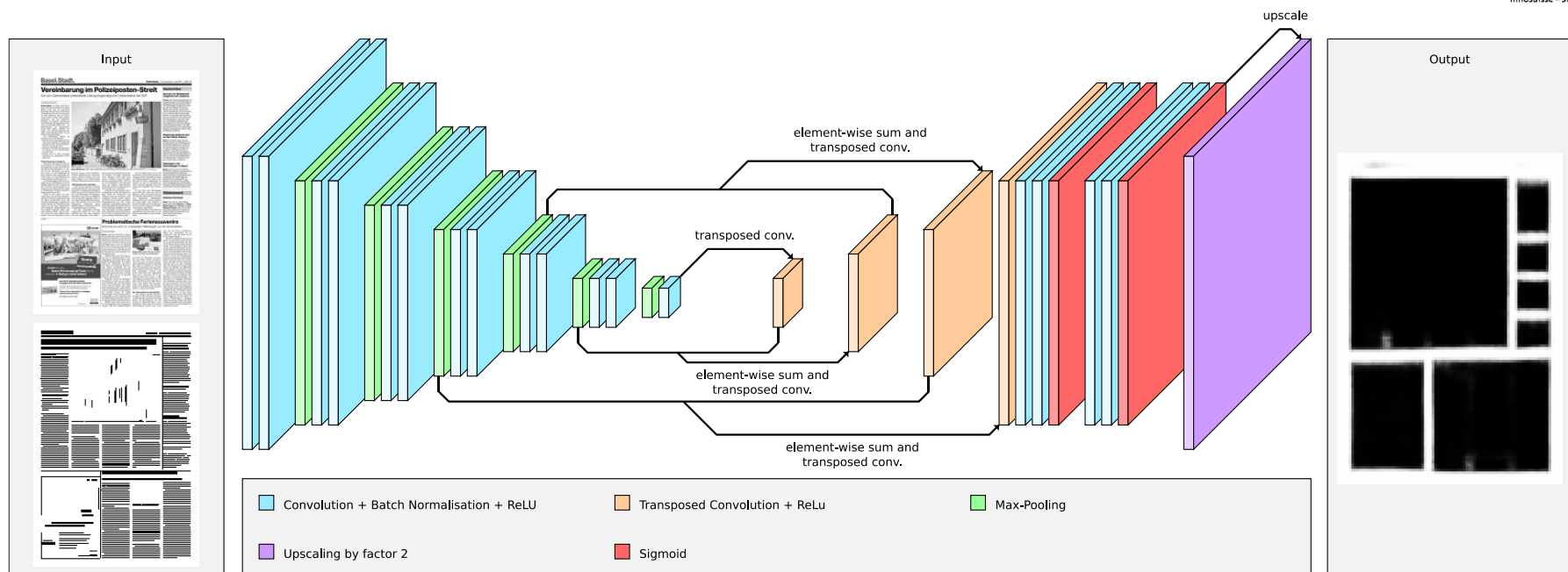
Sind Sie der nächste Lotto-König?

### Wochenpreis: 1x sieben Nächte für 2 Personen, inkl. HP, im \*\*\*\*Seehotel Pilatus Hergiswil im Wert von 3000 Franken!

9	4	6	2	5	1	7
3	1	6	9	8	4	1
4	6	2	8	2	4	4
6	7	0	6	8	8	4

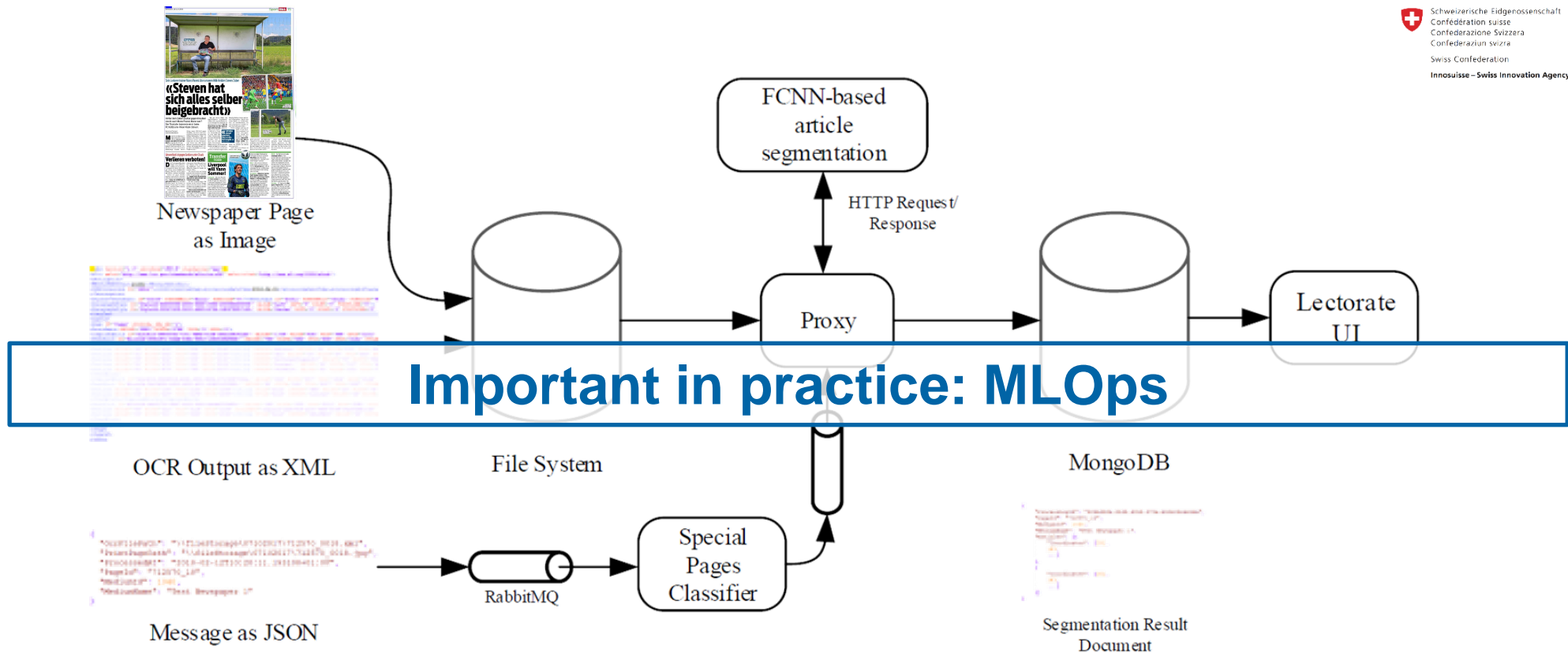
### GRÜBELMÄNNE

# 1. Print media monitoring – ML solution



Meier, Stadelmann, Stampfli, Arnold & Cieliebak (2017). «Fully Convolutional Neural Networks for Newspaper Article Segmentation». ICDAR'2017.  
Stadelmann, Tolkachev, Sick, Stampfli & Dürri (2018). «Beyond ImageNet - Deep Learning in Industrial Practice». In: Braschler et al., «Applied Data Science», Springer.

# 1. Print media monitoring – deployment



Stadelmann, Amirian, Arabaci, Arnold, Duivesteyn, Elezi, Geiger, Lörwald, Meier, Rombach & Tuggener (2018). «Deep Learning in the Wild». ANNPR'2018.

# 2. Music scanning

N 212

Die Forelle.  
Op. 28 No. 14. Scherz.  
Für eine Singstimme mit Begleitung des Pianoforte  
comp. aut. 1840

Schubert's Werk.  
FRANZ SCHUBERT.  
Erste Fassung.

Musik:  
Singsstimme:  
Pianoforte:



```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE score-partwise SYSTEM "http://www.musicxml.org/@/partwise.dtd" PUBLIC "-//Recordare/DTG MusicXML 2.0
Partwise/EN"
- <score-partwise>
- <identifications>
- <encoding>
- <software> MuseScore 1.3 </software>
- <encoding-date> 2014-12-16 </encoding-date>
- <encoding/>
- <source> http://musescore.com/score/502006 </source>
- <identification/>
- <defaults>
- <scaling>
- <millimeters> 7.056 </millimeters>
- <tenths> 40 </tenths>
- </scaling>
- </page-layout>
- <page-height> 1683.67 </page-height>
- <page-width> 1190.48 </page-width>
- <page-margins type="even">
- <left-margin> 56.6893 </left-margin>
- <right-margin> 56.6893 </right-margin>
- <top-margin> 56.6893 </top-margin>
- <bottom-margin> 113.379 </bottom-margin>
- </page-margins>
- <page-margins type="odd">
- <left-margin> 56.6893 </left-margin>
- <right-margin> 56.6893 </right-margin>
- <top-margin> 56.6893 </top-margin>
- <bottom-margin> 113.379 </bottom-margin>
- </page-margins>
- </page-layout>
- </defaults>
- <credit page="1">
- <credit words valign="top" justify="center" font-size="24" default-x="1626.98" default-y="595.238"> Die
Forelle </credit words>
- </credit>
- <credit page="1">
- <credit words valign="top" justify="right" font-size="12" default-y="1552.22" default-x="1133.79"> Franz
Schubert </credit words>
- </credit>
- <credit page="1">
- <credit words valign="bottom" justify="center" font-size="8" default-y="113.379" default-x="595.238"> Franz
Schubert, Die Forelle (Mollisande on http://www.Musescore.com) </credit words>
- </credit>
- <part-list>
- <score-part id="P1">
- <part name="Ténor" </part name>
- <part abbreviation="Ténor" </part abbreviation>
- <score-instrument id="P1-13">
- <instrument name="Ténor" </instrument name>
- </score-instrument>
- <midi-instrument id="P1-13">
- <midi-channel=1 </midi-channel>
- <midi-program=74 </midi-program>
- <volume=78.7402 </volume>
- <pan=0 </pan>
- </midi-instrument>
- </score-part>
- <part group type="start" number="1">
- <group-symbol=brace </group-symbol>
- </part group>
- <score-part id="P2">
- <part name="" </part name>
- <score-instrument id="P2-13">
- <instrument name="" </instrument name>
- </score-instrument>
- </part list>
```

Zürcher Hochschule für Angewandte Wissenschaften



Score Pad

Schweizerische Eidgenossenschaft  
Confédération suisse  
Confederazione Svizzera  
Confederaziun svizra  
Swiss Confederation  
Innosuisse – Swiss Innovation Agency



Die Forelle - Franz Schubert

$\text{♩} = 80$

Voice

Piano

Vo.

ei - nem Büch - lein hel - le, da schoß in fro - her Eil die lau - ni - sche Fo - re - le vor -



## 2. Music scanning – challenges & results

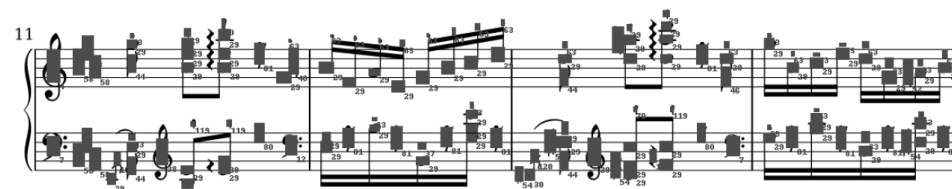


(a) accidentalSharp

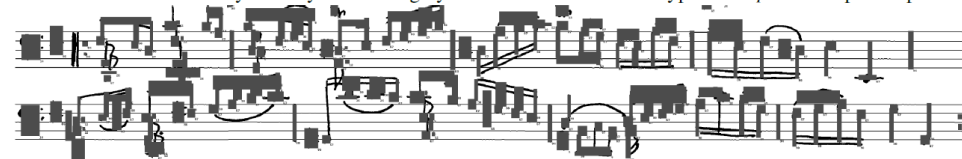
(b) keySharp

(c) augmentationDot

(d) articStaccatoAbove



a) Example result from *DeepScores* with detected bounding boxes as overlays. The tiny numbers are class labels from the dataset introduced with the overlay. This system is roughly one forth of the size of a typical *DeepScores* input we process at once.



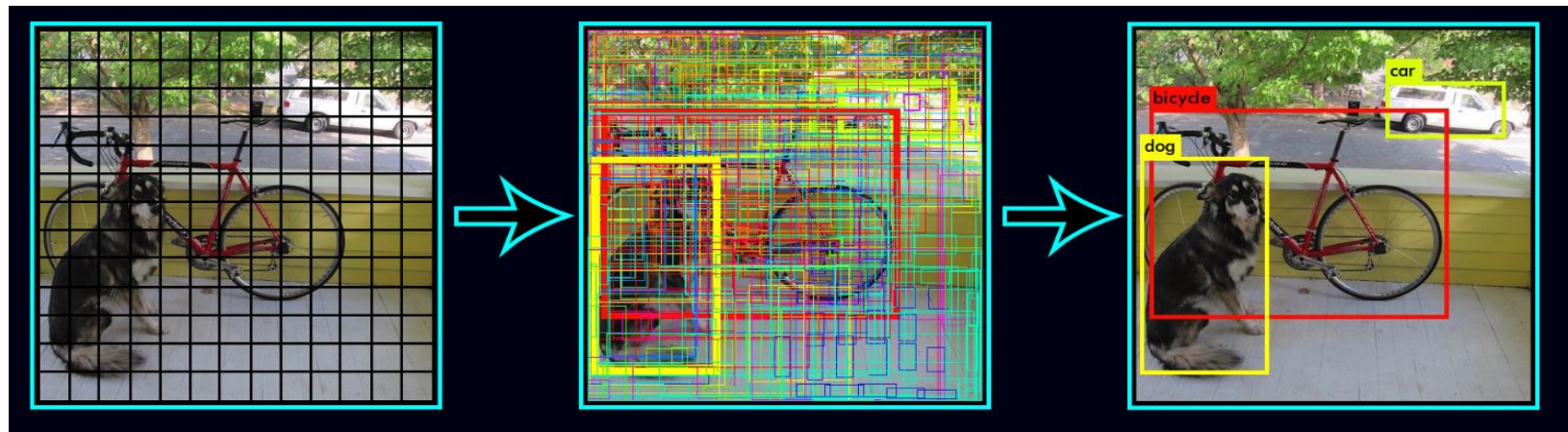
b) Example result from *MUSCIMA++* with detected bounding boxes and class labels as overlays. This system is roughly one half of the size of a typical processed *MUSCIMA++* input. The images are random picks amongst inputs with many symbols.

Tuggener, Elezi, Schmidhuber, Pelillo & Stadelmann (2018). «*DeepScores – A Dataset for Segmentation, Detection and Classification of Tiny Objects*». ICPR'2018.  
 Tuggener, Satyawan, Pacha, Schmidhuber & Stadelmann (2020). «*The DeepScoresV2 Dataset and Benchmark for Music Object Detection*». ICPR'2020.  
 Tuggener, Elezi, Schmidhuber & Stadelmann (2018). «*Deep Watershed Detector for Music Object Recognition*». ISMIR'2018.

# Music scanning – methodology (differentiation)


## OMR vs state of the art object detectors

### YOLO/SSD-type detectors



Source: <https://pjreddie.com/darknet/yolov2/> (11.09.2018)

Score Pad

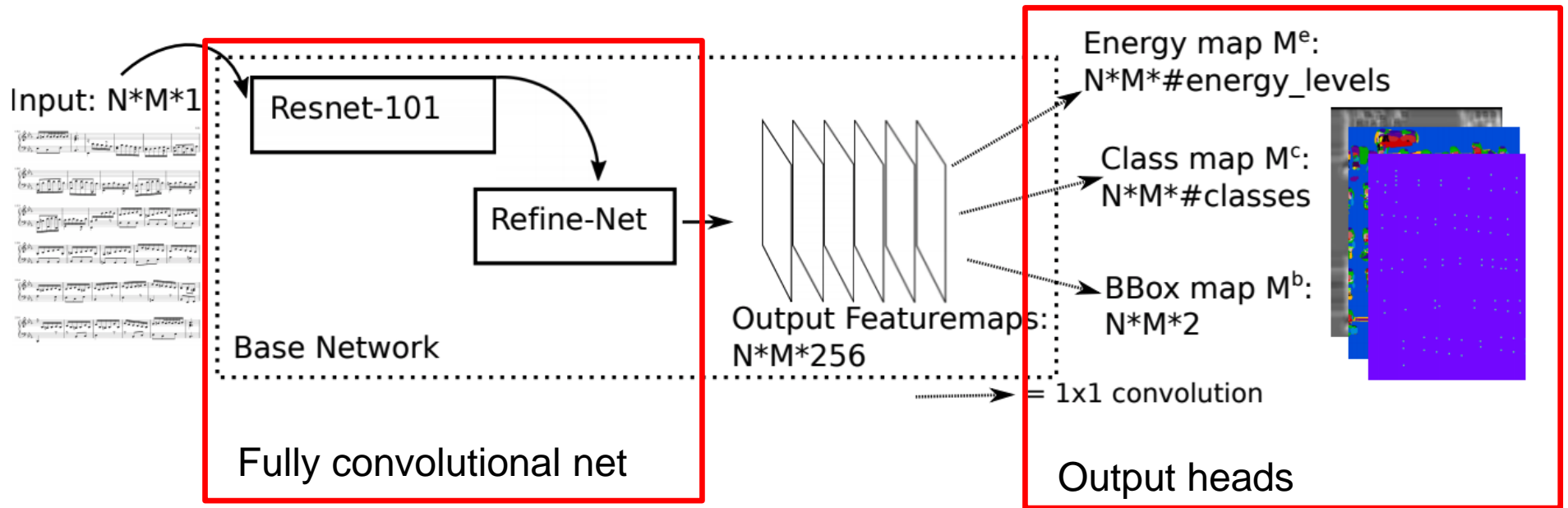
 Schweizerische Eidgenossenschaft  
Confédération suisse  
Confederazione Svizzera  
Confederaziun svizra  
Swiss Confederation  
Innosuisse – Swiss Innovation Agency

### R-CNN

- Two-step proposal and refinement scheme
- Very large number of proposals at high resolution needed

# 2. Music scanning – methodology (ours)

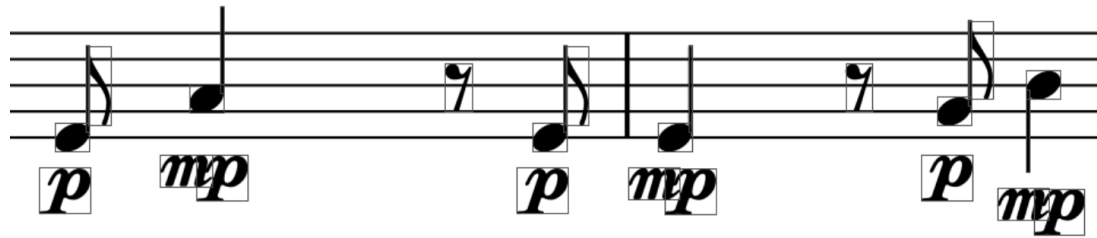
## The deep watershed detector



# 2. Music scanning – methodology (details)

## Output heads of the deep watershed detector

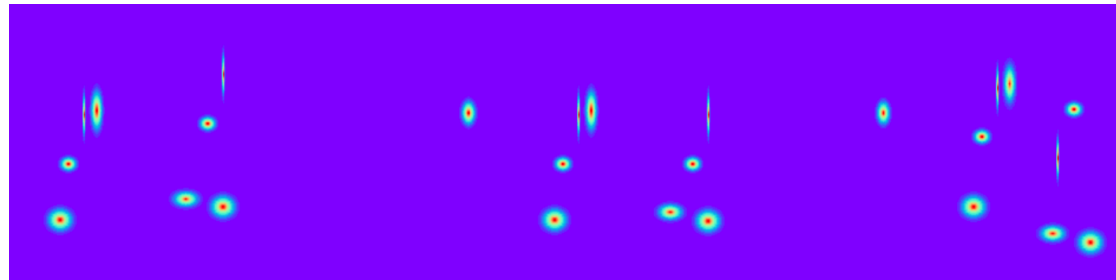
Input



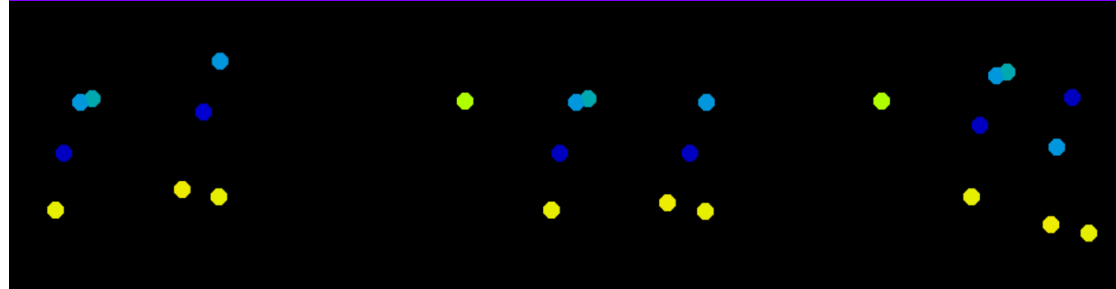
Score Pad

Schweizerische Eidgenossenschaft  
Confédération suisse  
Confederazione Svizzera  
Confederaziun svizra  
Swiss Confederation  
Innosuisse – Swiss Innovation Agency

Energy

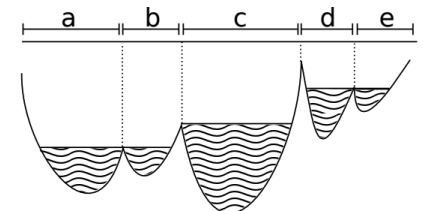
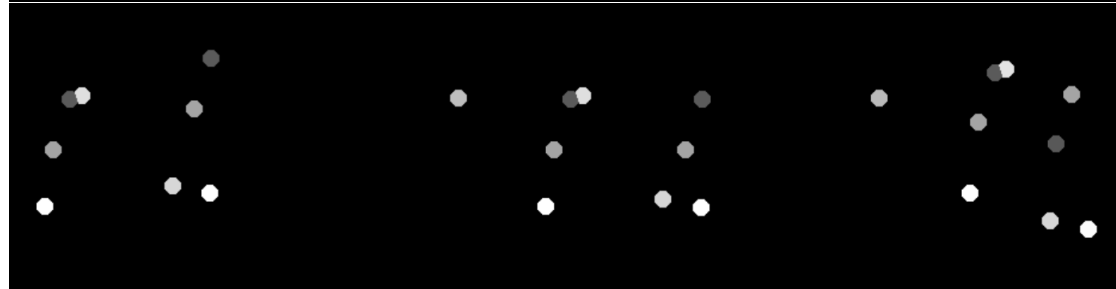


Class

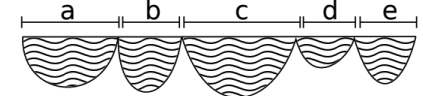


Size

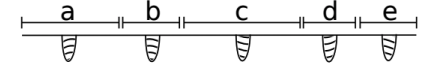
2D projection



a) One-dimensional energy function of five classes without any structural constraints.



b) Energy function for the same five classes with fixed boundary energy.



c) Energy function for the same five classes this time with small energy markers at the class centers.

## 2. Music scanning – recent work

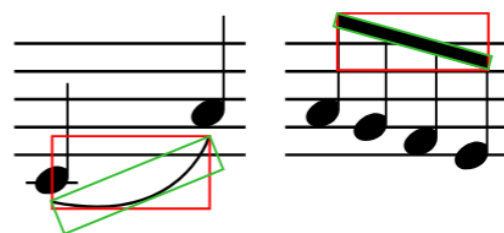
1. Extension of the detection alphabet to 135 classes (updated dataset released)



Score Pad

Schweizerische Eidgenossenschaft  
Confédération suisse  
Confederazione Svizzera  
Confederaziun Svizra  
Swiss Confederation  
Innosuisse – Swiss Innovation Agency

2. Natural incorporation of rotation



## 2. Music scanning – future work

### Dealing with real world noise

Synthetic quality + labels:  
perfect quality

Real world quality: print/scan  
artifacts, wrinkles, dirt, ...

Score Pad



Schweizerische Eidgenossenschaft  
Confédération suisse  
Confederazione Svizzera  
Confederaziun svizra  
Swiss Confederation

Innosuisse – Swiss Innovation Agency

Model training

**Data distribution shift**

Model deployment

Remedy:

Use semi-supervised learning to  
model distribution change and/or  
disentangle latent signals

### 3. Outlook – lessons learned

Data is key

- Many real-world projects miss the required **quantity & quality** of data  
→ even though «big data» is not needed
- **Class imbalance** needs careful dealing  
→ special loss, resampling (also in unorthodox ways)
- **Unsupervised** methods need to be used creatively
- Users & label providers need to be **trained**

**Prerequisite: stable data & label acquisition pipeline**

**Learning from (raw) data is powerful, yet one is fully dependent on what is in that data**

Robustness is important

- **Training processes** can be tricky  
→ give hints via a unique loss, proper preprocessing and pretraining

**Sufficient condition: lots of tuning**

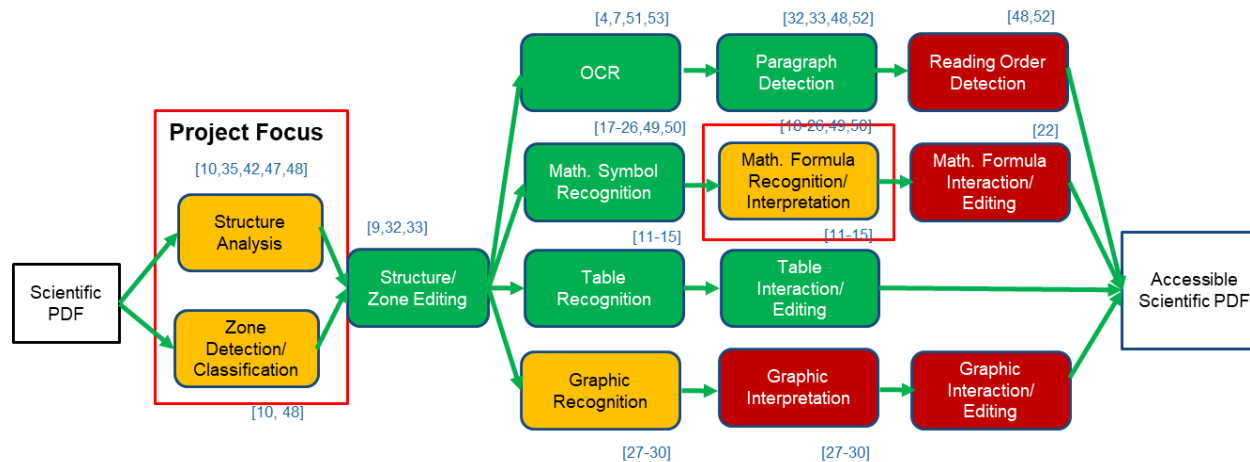
# 3. Outlook – accessible documents for all

The screenshot shows the PAVE application interface. The top navigation bar includes the PAVE logo and links for ZHAW, CONTACT, HELP/FAQ, and LANGUAGE. Below this, the document title 'sample1\_result.pdf' is displayed. A menu with 'TASKS', 'PROPERTIES', 'ISSUE DETAILS', and 'READING ORDER' is visible. The 'READING ORDER' tab is active, showing a list of document sections with checkboxes and page numbers. One section, 'Paragraph  $\sigma \times \sigma(X(\sigma)(\sigma)(\sigma)P=vE, P=vEfi9k+1ii$ ', is highlighted with a yellow circle. The main content area displays 'Page 1' of the document, which includes a title 'MASAKAZU SUZUKI', a proof section, and mathematical diagrams and equations. The diagram shows a sum of terms with subscripts, and the equations involve summations over  $j$  and  $i$ .

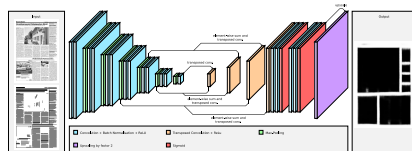


### 3. Outlook – research approach

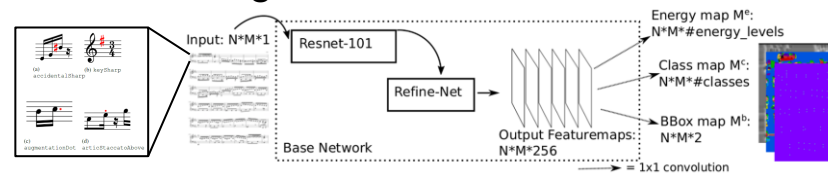
- Goal is a **semi-automatic process** to make scientific documents accessible
- Focus on **structure analysis & mathematical formula recognition**



- **Build on previous experience & methods** for newspapers and OMR structure recognition



formula recognition



# Contact

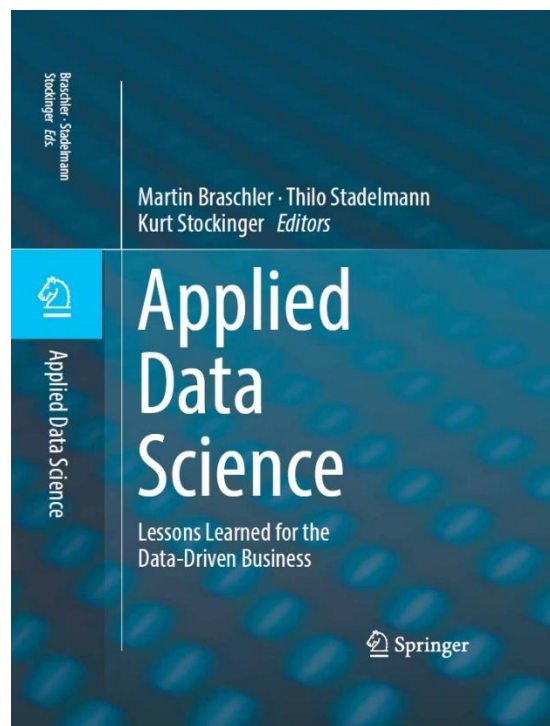
- Document analysis is a **very fruitful use case** for Deep Learning (for archives + R&D)
- **Latest research is applied** and deployed in «normal» organizations (e.g., libraries)
- It does not need big-, but some **data (effort usually underestimated)**
- DL/RL **training** for new use cases **can be tricky** (→ needs experience)

## About me:

- Prof. AI/ML, scientific director ZHAW digital
- Email: [stdm@zhaw.ch](mailto:stdm@zhaw.ch)
- Phone: +41 58 934 72 08
- Web: <https://stdm.github.io/>
- Twitter: [@thilo\\_on\\_data](https://twitter.com/thilo_on_data)
- LinkedIn: [thilo-stadelmann](https://www.linkedin.com/in/thilo-stadelmann)

## Further contacts:

- [datalab@zhaw.ch](mailto:datalab@zhaw.ch), [info.office@data-innovation.org](mailto:info.office@data-innovation.org), [office-switzerland@claire-ai.org](mailto:office-switzerland@claire-ai.org)





# APPENDIX

# Team AI/ML: Overview (cp. <https://stdm.github.io/research/>)

## ZHAW School of Engineering, Winterthur, Switzerland [2]



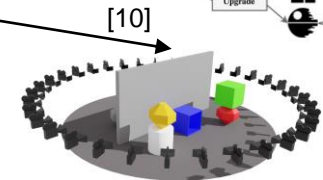
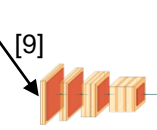
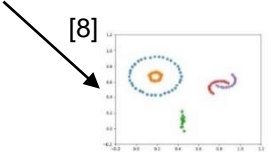
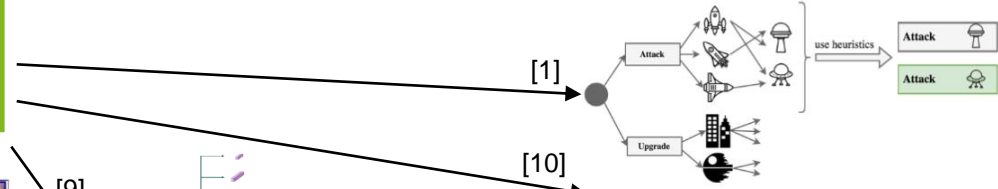
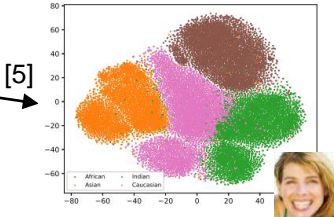
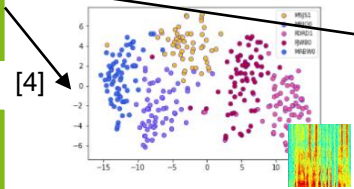
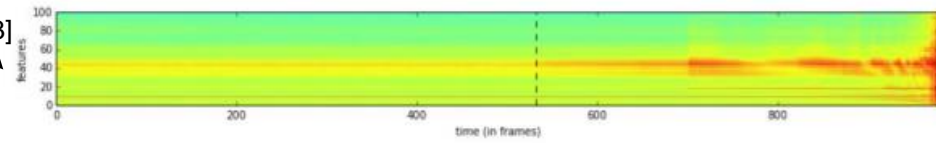
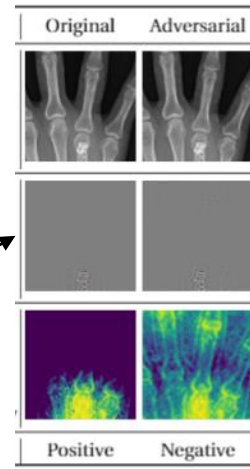
### Machine learning-based Pattern Recognition

Robust applications

Biometrics

Document Analysis

Learning to act

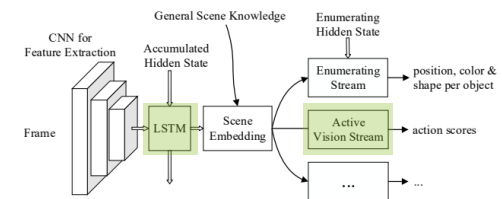
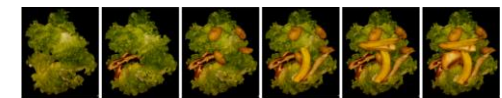
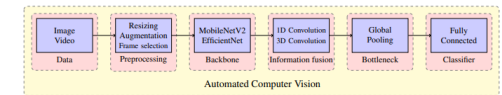
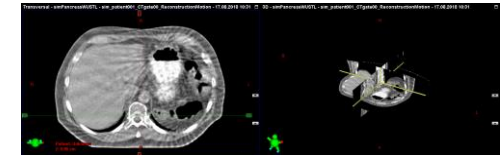


# References for overview

1. Thilo Stadelmann, Mohammadreza Amirian, Ismail Arabaci, Marek Arnold, Gilbert François Duivesteijn, Ismail Elezi, Melanie Geiger, Stefan Lörwald, Benjamin Bruno Meier, Katharina Rombach, and Lukas Tuggener. **“Deep Learning in the Wild”**. In: Proceedings of the 8th IAPR TC 3 Workshop on Artificial Neural Networks for Pattern Recognition (**ANNPR’18**), Springer, LNAI 11081, pp. 17-38, Siena, Italy, September 19-21, 2018.
2. Mohammadreza Amirian, Friedhelm Schwenker, and Thilo Stadelmann. **“Trace and Detect Adversarial Attacks on CNNs using Feature Response Maps”**. In: Proceedings of the 8th IAPR TC 3 Workshop on Artificial Neural Networks for Pattern Recognition (**ANNPR’18**), Springer, LNAI 11081, pp. 346-358, Siena, Italy, September 19-21, 2018.
3. Thilo Stadelmann, Vasily Tolkahev, Beate Sick, Jan Stampfli, and Oliver Dürr. **“Beyond ImageNet - Deep Learning in Industrial Practice”**. In: Martin Braschler, Thilo Stadelmann, and Kurt Stockinger (Editors). **“Applied Data Science - Lessons Learned for the Data-Driven Business”**. Springer, 2019.
4. Thilo Stadelmann, Sebastian Glinski-Haefeli, Patrick Gerber, and Oliver Dürr. **“Capturing Suprasegmental Features of a Voice with RNNs for Improved Speaker Clustering”**. In: Proceedings of the 8th IAPR TC 3 Workshop on Artificial Neural Networks for Pattern Recognition (**ANNPR’18**), Springer, LNAI 11081, pp. 333-345, Siena, Italy, September 19-21, 2018.
5. Stefan Glüge, Mohammadreza Amirian, Dandolo Flumini, and Thilo Stadelmann. **“How (Not) to Measure Bias in Face Recognition Networks”**. In: Proceedings of the 9th IAPR TC 3 Workshop on Artificial Neural Networks for Pattern Recognition (**ANNPR’20**), Springer, LNAI, Winterthur, Switzerland, September 02-04, 2020.
6. Lukas Tuggener, Yvan Putra Satyawan, Alexander Pacha, Jürgen Schmidhuber, and Thilo Stadelmann. **“The DeepScoresV2 Dataset and Benchmark for Music Object Detection”**. In: Proceedings of the 25th International Conference on Pattern Recognition (**ICPR’20**), IAPR, Milan, Italy, January 10-15 (online), 2021.
7. Benjamin Meier, Thilo Stadelmann, Jan Stampfli, Marek Arnold, and Mark Cieliebak. **“Fully convolutional neural networks for newspaper article segmentation”**. In: Proceedings of the 14th IAPR International Conference on Document Analysis and Recognition (**ICDAR’17**). 14th IAPR International Conference on Document Analysis and Recognition (ICDAR), Kyoto Japan, November 13-15, 2017. Kyoto, Japan: CPS.
8. Benjamin Bruno Meier, Ismail Elezi, Mohammadreza Amirian, Oliver Dürr, and Thilo Stadelmann. **“Learning Neural Models for End-to-End Clustering”**. In: Proceedings of the 8th IAPR TC 3 Workshop on Artificial Neural Networks for Pattern Recognition (**ANNPR’18**), Springer, LNAI 11081, pp. 126-138, Siena, Italy, September 19-21, 2018.
9. Lukas Tuggener, Mohammadreza Amirian, Fernando Benites, Pius von Däniken, Prakhar Gupta, Frank-Peter Schilling, and Thilo Stadelmann. **“Design Patterns for Resource-Constrained Automated Deep-Learning Methods”**. AI section “Intelligent Systems: Theory and Applications” 1(4):510-538, MDPI, Basel, Switzerland, November 06, 2020.
10. Dano Roost, Ralph Meier, Giovanni Toffetti Carughi, and Thilo Stadelmann. **“Combining Reinforcement Learning with Supervised Deep Learning for Neural Active Scene Understanding”**. In: Proceedings of the Active Vision and Perception in Human(-Robot) Collaboration Workshop at IEEE RO-MAN 2020 (**AVHRC’20**), online, August 31, 2020.

# Outlook: Late-breaking results

- Medical image analysis: learning to reduce motion artifacts in 3D CT scans
- Learning an artificial communication language for multi-agent reinforcement learning in logistics (notable rank in Flatland 2019 competition, best poster award [1])
- Automated deep learning (top rank in AutoDL 2020 challenge [2])
- Learning to segment and classify food waste in professional kitchens under adversarial conditions [4]
- Improving robotic vision through active vision and combined supervised and reinforcement learning (Dr. Waldemar Jucker Award 2020 [3])



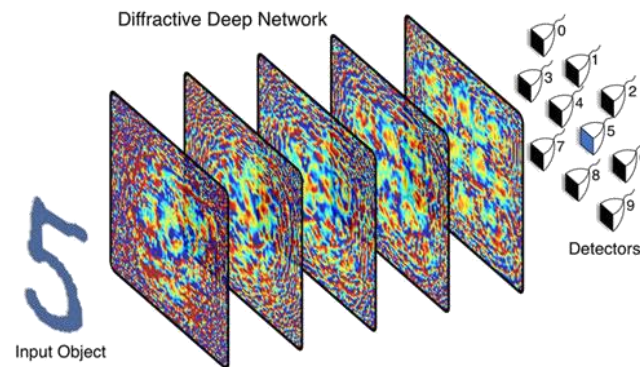
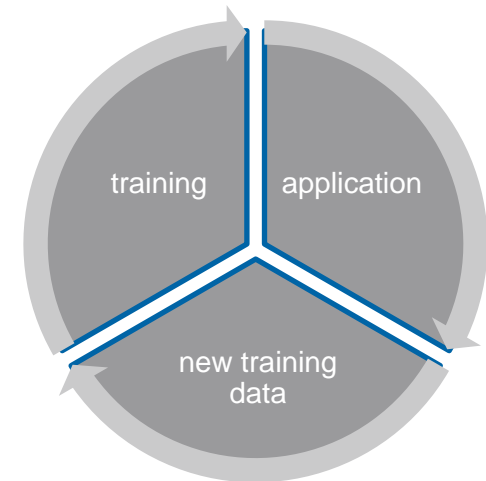
[1] Roost, Meier, Huschauer, Nygren, Egli, Weiler & Stadelmann (2020). «Improving Sample Efficiency and Multi-Agent Communication in RL-based Train Rescheduling». SDS'2020.  
 [2] Tuggener, Amirian, Benites, von Däniken, Gupta, Schilling & Stadelmann (2020). «Design Patterns for Resource Constrained Automated Deep Learning Methods». AI 1(4) 510-538.  
 [3] Roost, Meier, Toffetti Carughi & Stadelmann (2020). «Combining Reinforcement Learning with Supervised Deep Learning for Neural Active Scene Understanding». AVHRC 2020.  
 [4] Simmler, Sager, Andermatr, Chavarriaga, Schilling, Rosenthal & Stadelmann (2021). «Noisy labels, missing labels: A survey of un-, weakly-, and semi-supervised learning methods for industrial vision applications». SDS 2021.

# 1. Lessons learned 1/2



## Deployment

- Should include **continuous learning**
- Needs to take care of **processing speed / efficiency**



Symbolic image: a CNN in (optical) hardware (Lin et al., 2018).

Lin, Rivenson, Yardimci, Veli, Luo, Jarrahi & Oczan (2018). «All-optical machine learning using diffractive deep neural networks». Science, 26. Jul 2018.

## 2. Music scanning – methodology (further details)

### Reconstructing labeled bounding boxes

#### Finding object instances:

1. Set all values of the energy prediction below a certain threshold to zero (watershed step).
2. Perform a connected component analysis on this output --> the center of mass for every connected component is defined to be the position of one detected object.

#### Predicting the class:

3. For each connected component look up the corresponding class predictions for of all its pixels and use the majority vote.

#### Predicting object size:

4. For each connected component look up the corresponding size predictions for of all its pixels and compute the mean.



## Lessons learned 2/2



### Loss shaping

- Usually necessary to **enable learning** of very complex target functions

*“Initially, the training was **unstable** [...] if directly trained on the **combined weighted loss**. Therefore, we now **train** [...] on each of the **three tasks separately**.*

*We further observed that while the network gets trained on the bounding box prediction and classification, the energy level predictions get worse. To avoid this, the network is **fine-tuned only for the energy level loss** [...]. Finally, the network is retrained on the combined task [...] for a few thousand iterations [...].”*

- This includes **encoding expert knowledge** manually into the model architecture or training setup

*“The **size of the anomaly** in classifying balloon catheters as good or bad is **quite decisive**. Thus, rescaling the training images is not allowed, and we used a fixed size window around the center of each defect to extract the training images.”*

Stadelmann, Amirian, Arabaci, Arnold, Duivesteyn, Elezi, Geiger, Lörwald, Meier, Rombach & Tuggener (2018). «Deep Learning in the Wild». ANNPR'2018.