

Deep Learning-based Pattern Recognition in Business

40. Berner Architektentreffen , June 29, 2018

Thilo Stadelmann



Swiss Alliance for
Data-Intensive Services

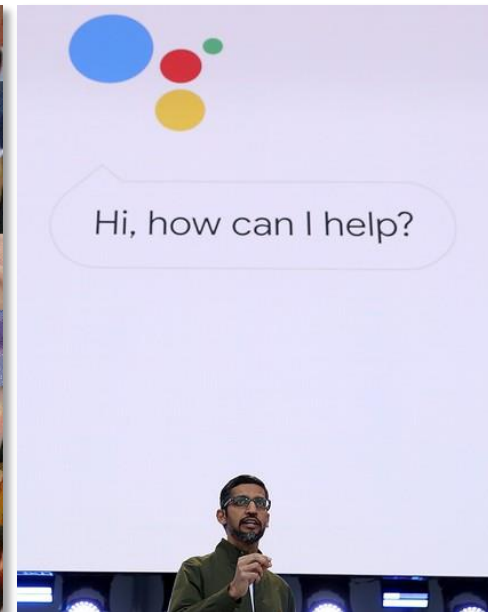
swiss group for artificial intelligence
and cognitive science



data lab

www.zhaw.ch/data lab

Why?



Why? – deep learning in a nutshell

Classical
pattern recognition

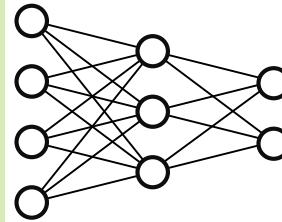


Hand-crafted features
(SIFT, SURF, LBP, HOG, etc.)

(0.2, 0.4, ...)

(0.4, 0.3, ...)

Classifikation
(SVM, neuronal net, etc.)



Container ship

Tiger

...

Why? – deep learning in a nutshell

Classical
pattern recognition

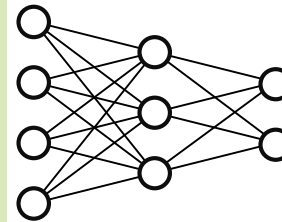


Hand-crafted features
(SIFT, SURF, LBP, HOG, etc.)

(0.2, 0.4, ...)

(0.4, 0.3, ...)

Classifikation
(SVM, neuronal net, etc.)



Container ship

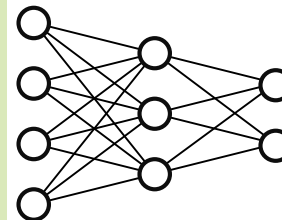
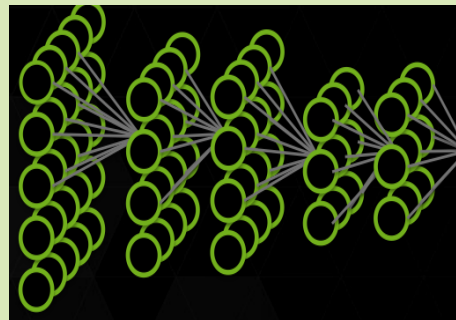
Tiger

...

Convolutional
neural network



Learns salient features from
pure pixels

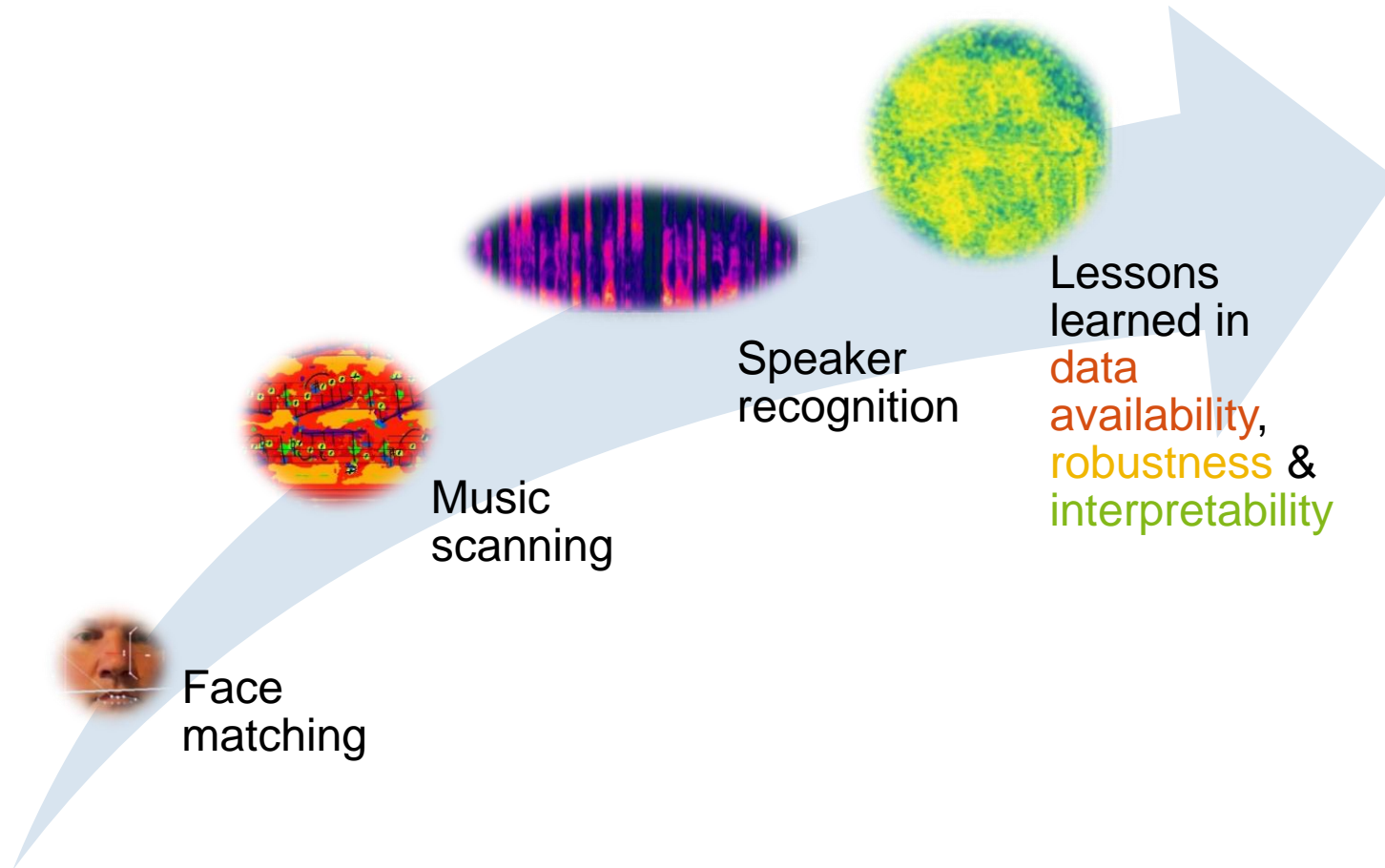


Container ship

Tiger

...


Agenda



Face matching




 **DEEPIMPACT**

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

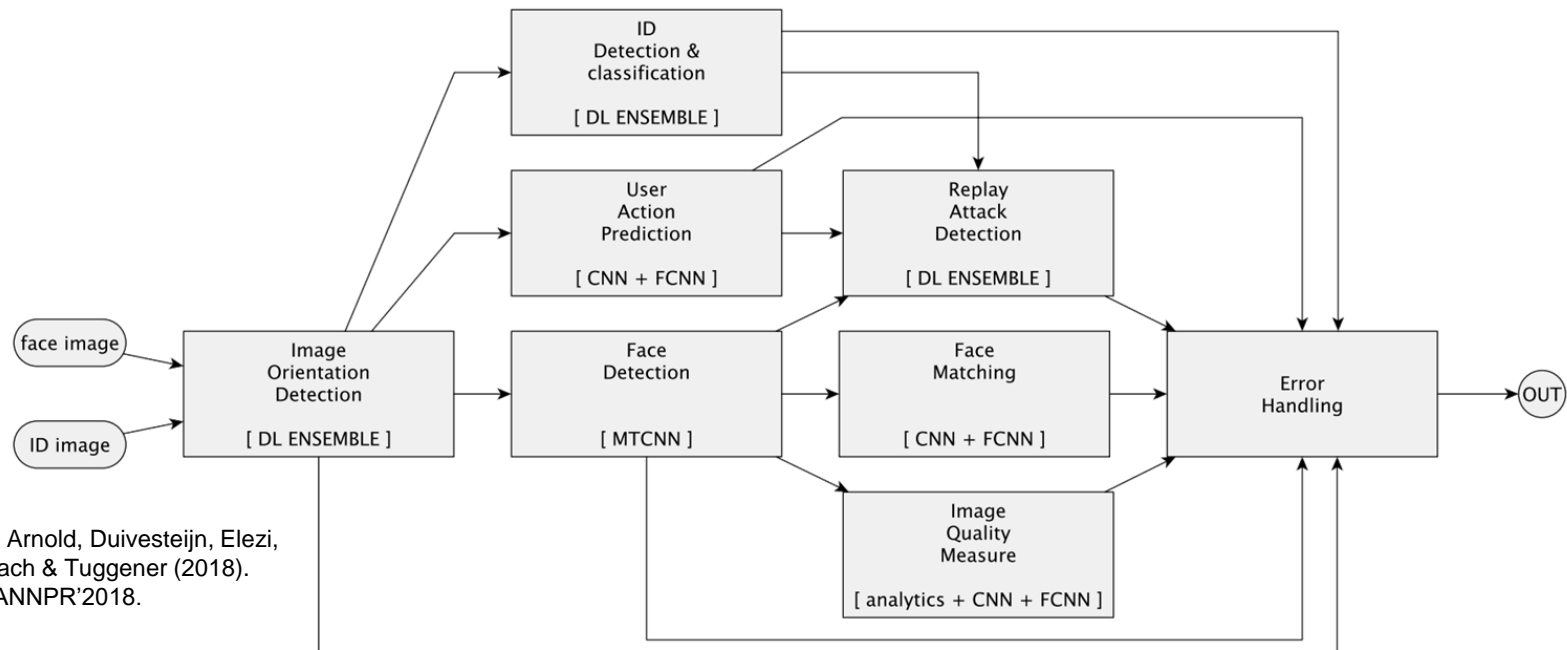
Face matching



 **DEEPIMPACT**

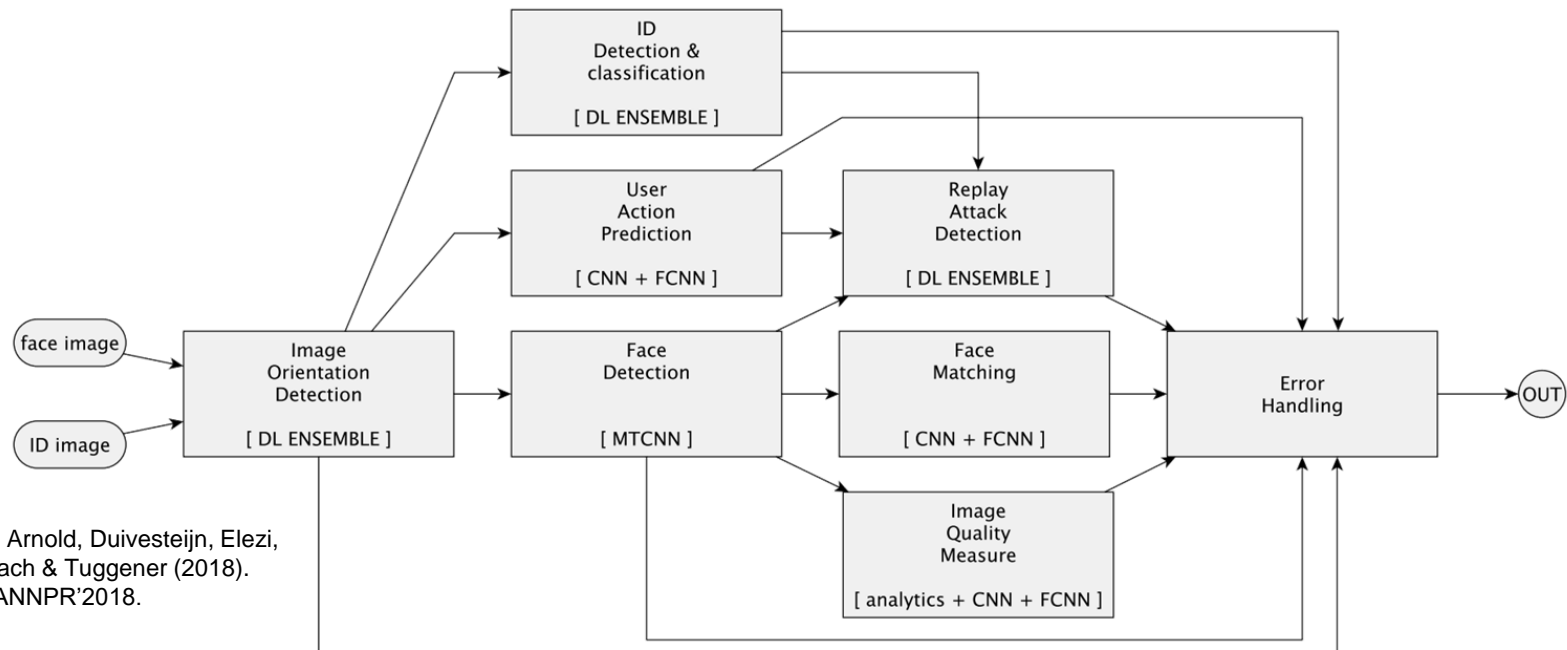
 Schweizerische Eidgenossenschaft
Confédération suisse
Confederazione Svizzera
Confederaziun svizra
Swiss Confederation
Innosuisse – Swiss Innovation Agency

Face matching – challenges & solutions



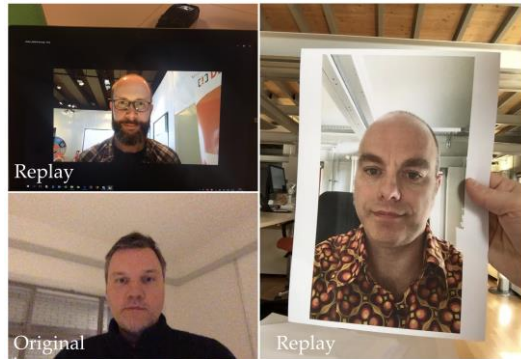
Stadelmann, Amirian, Arabaci, Arnold, Duivesteijn, Elezi, Geiger, Lörwald, Meier, Rombach & Tuggener (2018). «Deep Learning in the Wild». ANNPR'2018.

Face matching – challenges & solutions



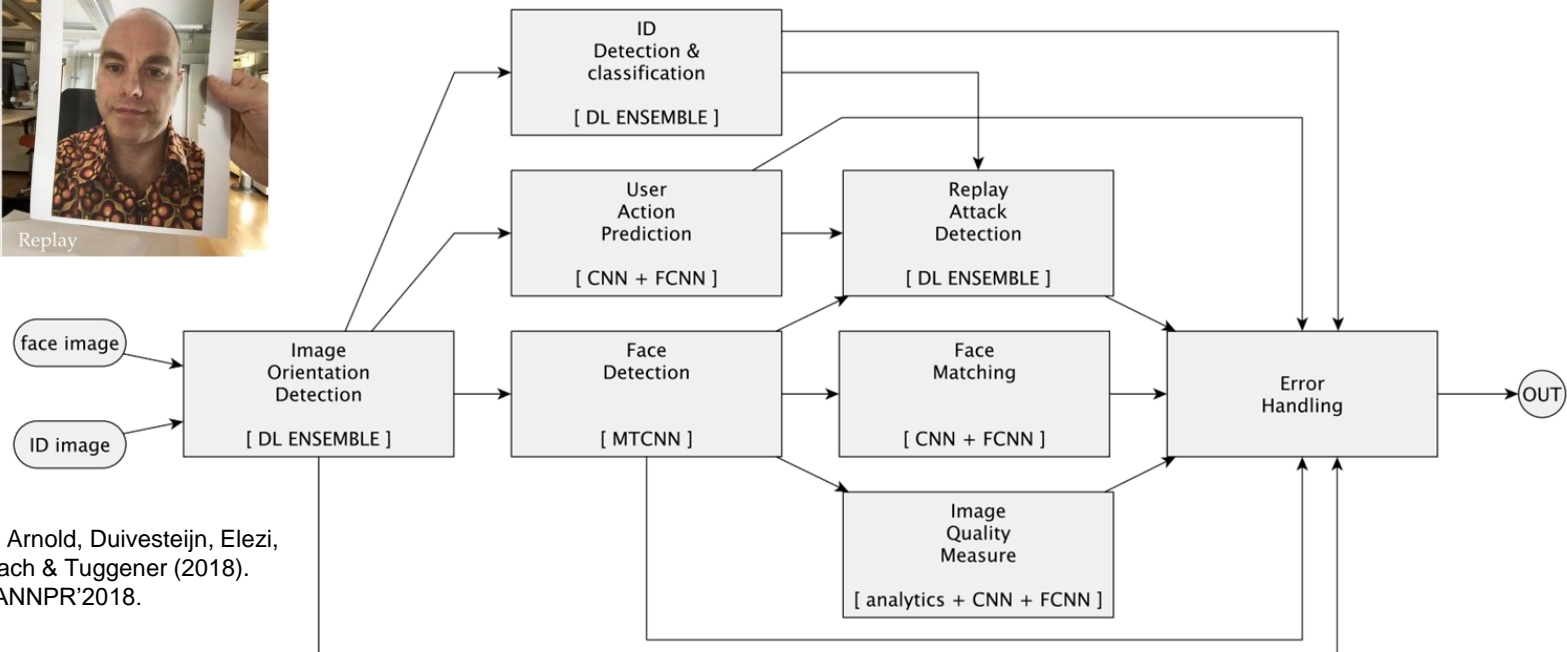
Stadelmann, Amirian, Arabaci, Arnold, Duivesteijn, Elezi, Geiger, Lörwald, Meier, Rombach & Tuggener (2018). «Deep Learning in the Wild». ANNPR'2018.

Face matching – challenges & solutions



[!] DEEPIIMPACT

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



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

Music scanning

N 212.
 Die Forelle.
 Op. 101 No. 35.
 Für eine Singstimme mit Begleitung des Pianoforte.
 Schuber's Werk.
 Franz Schubert.
 Erste Fassung.
 N° 212

Musik:
 Singstimme:
 Pianoforte:



```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE score-partwise SYSTEM "http://www.musescore.org/@id/partwise.dtd" PUBLIC "-//Recordare/DTG MusicXML 2.0 Partwise/EN"
- <score-partwise>
- <identification>
- <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>
- <cenths> 40 </cenths>
- </scaling>
- </page-layout>
- <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>
- <credit words valign="top" justify="center" font-size="24" default-y="1626.98" default-x="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>
- <score-instrument id="P2-13">
- <instrument-name>
```



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 -

Music scanning – challenges & solutions



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

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

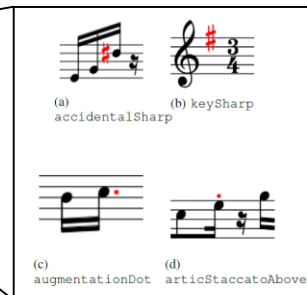
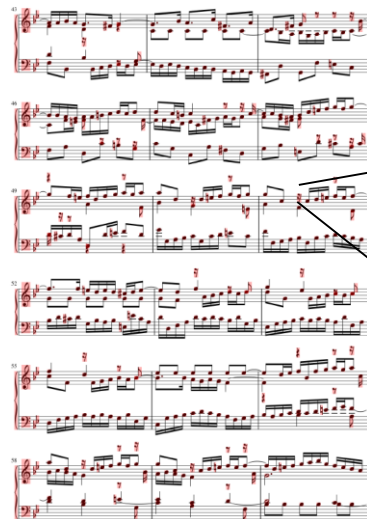
Music scanning – challenges & solutions

The image shows a musical score with several staves. A callout box highlights four specific annotations:

- (a) accidentalSharp
- (b) keySharp
- (c) augmentationDot
- (d) articStaccatoAbove

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

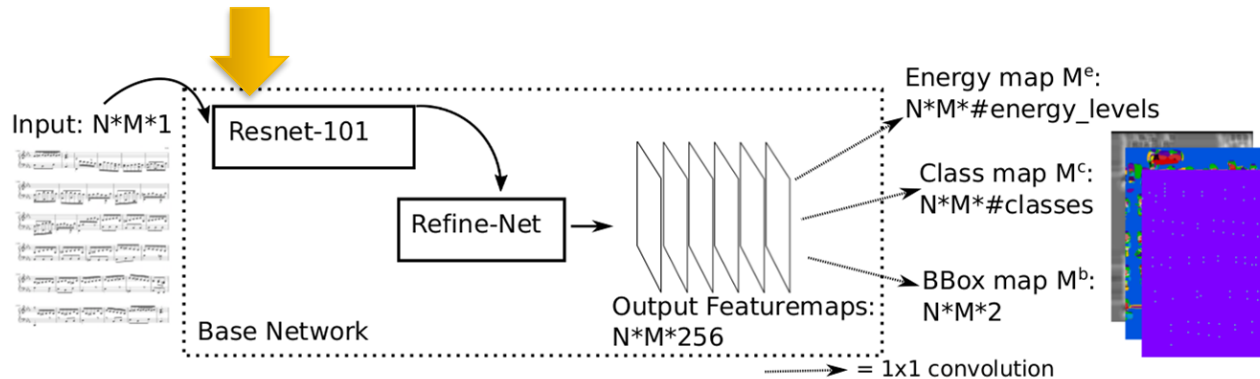
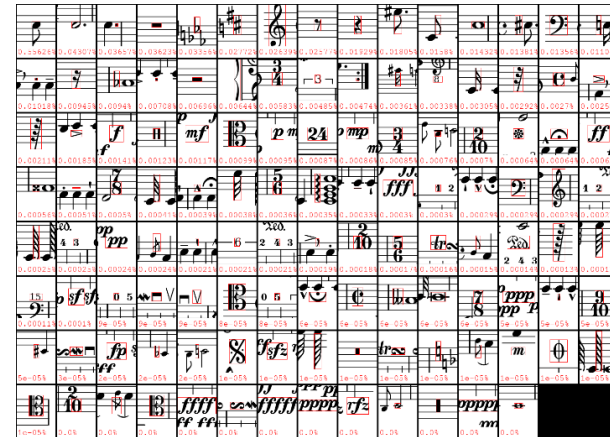
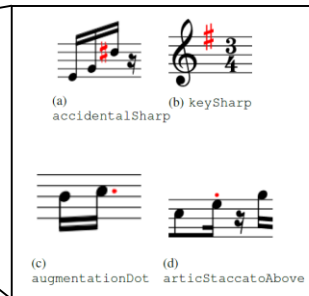
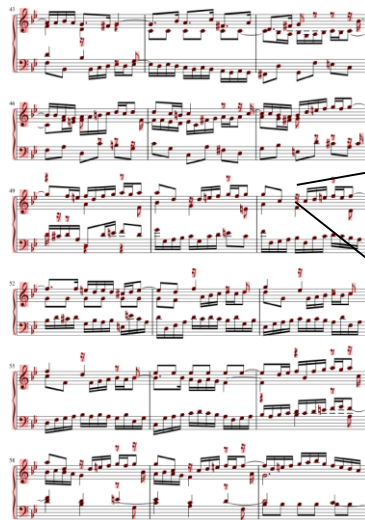
Music scanning – challenges & solutions



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

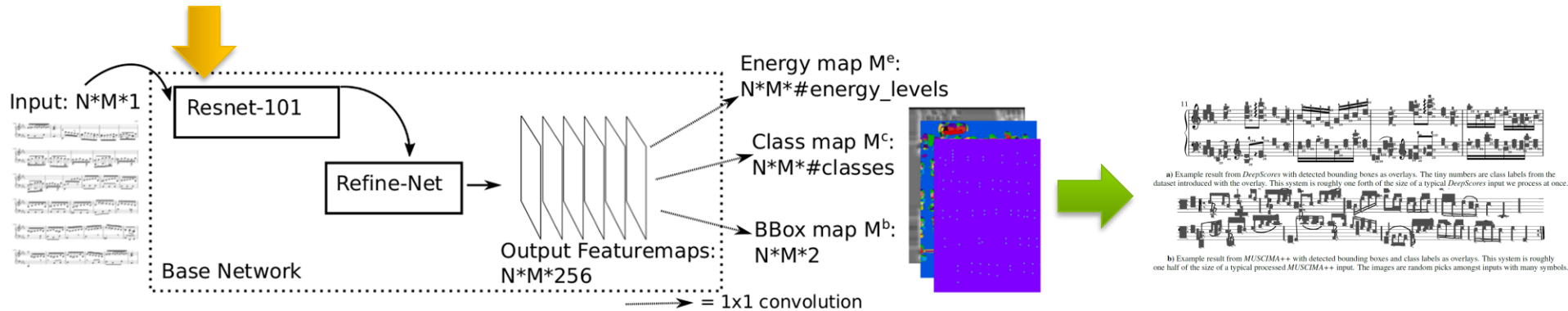
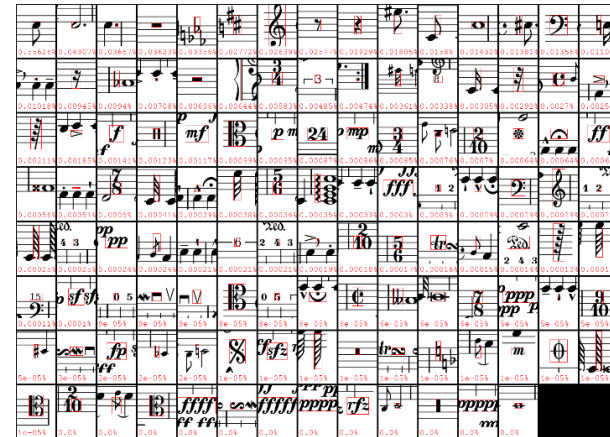
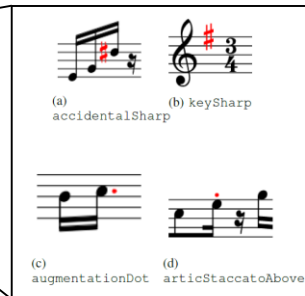
Tuggener, Elezi, Schmidhuber, Pelillo & Stadelmann (2018). «DeepScores – A Dataset for Segmentation, Detection and Classification of Tiny Objects». ICPR'2018.
Tuggener, Elezi, Schmidhuber & Stadelmann (2018). «Deep Watershed Detector for Music Object Recognition». ISMIR'2018.

Music scanning – challenges & solutions



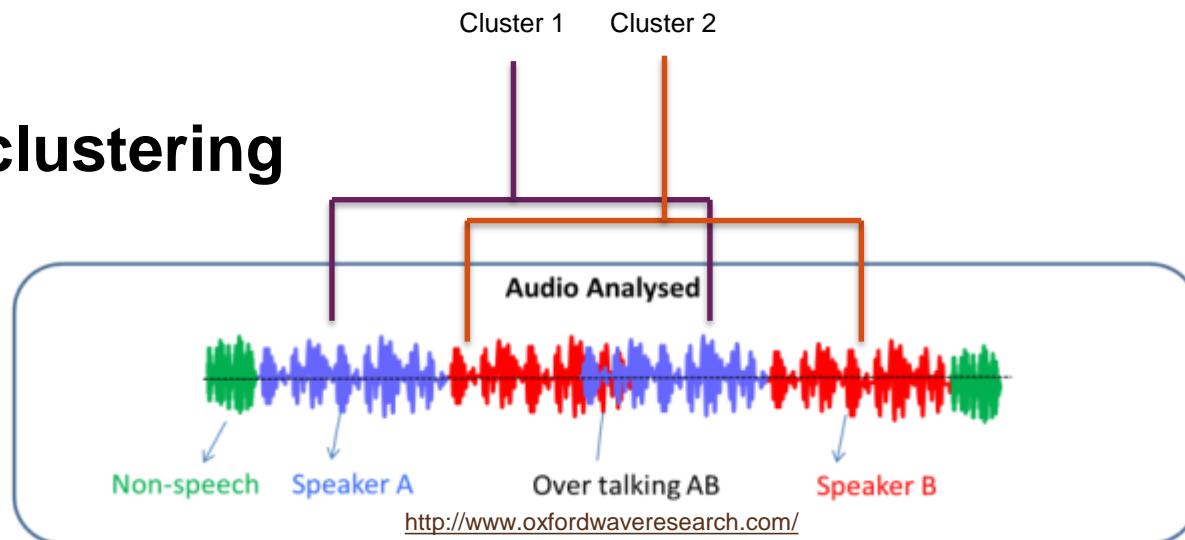
Tuggener, Elezi, Schmidhuber, Pelillo & Stadelmann (2018). «DeepScores – A Dataset for Segmentation, Detection and Classification of Tiny Objects». ICPR'2018.
Tuggener, Elezi, Schmidhuber & Stadelmann (2018). «Deep Watershed Detector for Music Object Recognition». ISMIR'2018.

Music scanning – challenges & solutions



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

Speaker clustering

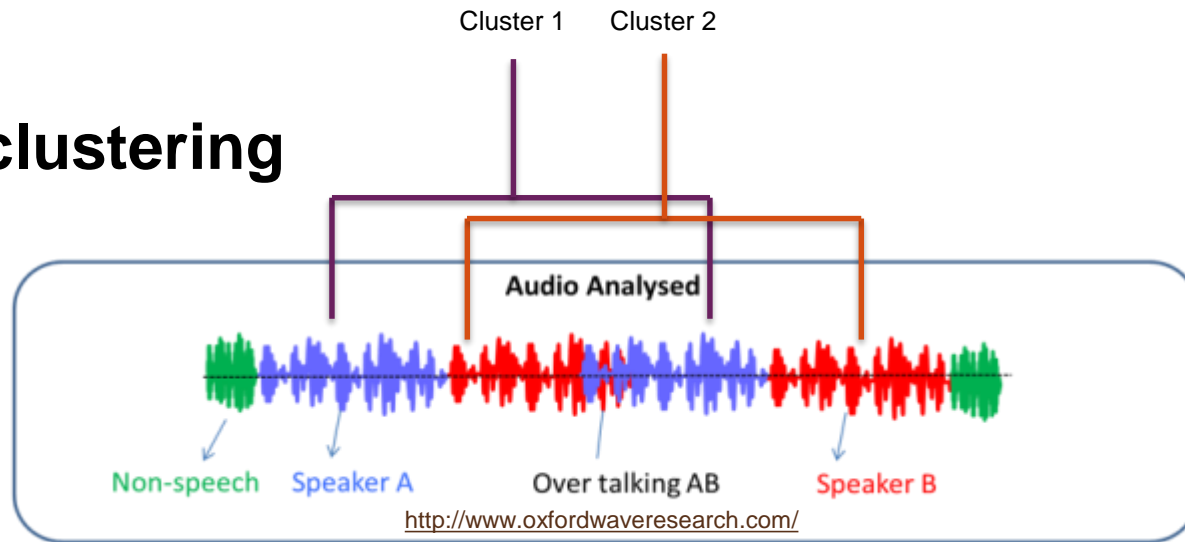


For the 630 training utterances, GMMs with 32 mixtures are built a priori, then an identification experiment is run for the 630 test utterances. It yields a satisfactory 0.5% closed set identification error.

[34]. Evaluations typically concentrate on data sets built from broadcast news/shows and meeting recordings, where diarization error rates ranging from 8% to 24% are reported [28][34][45]. These results are confirmed by more recent

The hypothesis of this paper is: the techniques originally developed for speaker verification and identification are not suitable for speaker clustering, taking into account the escalated difficulty of the latter task. However, the processing chain for speaker clustering is quite large – there are many potential areas for improvement. The question is: *where should improvements be made to improve the final result?*

Speaker clustering



For the 630 training utterances, GMMs with 32 mixtures are built a priori, then an identification experiment is run for the 630 test utterances. It yields a satisfactory 0.5% closed set identification error.

[34]. Evaluations typically concentrate on data sets built from broadcast news/shows and meeting recordings, where diarization error rates ranging from 8% to 24% are reported [28][34][45]. These results are confirmed by more recent

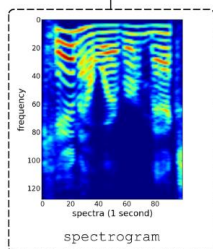
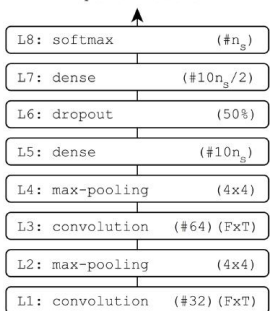
The hypothesis of this paper is: the techniques originally developed for speaker verification and identification are not suitable for speaker clustering, taking into account the escalated difficulty of the latter task. However, the processing chain for speaker clustering is quite large – there are many potential areas for improvement. The question is: *where should improvements be made to improve the final result?*

The interpretation of our results has shown that it is the stage of modeling that bears the highest potential: the inclusion of temporal context information among feature vectors is what is crucially missing there. Furthermore, the inclusion

context vector. This corresponds to a syllable length of 130 ms and is found to best capture speaker specific sounds in informal listening experiments over a range of 32–496 ms (in intervals of 16 ms). Our context vector step is one orig-

Speaker clustering – exploiting time information

CNN (MLSP'16)



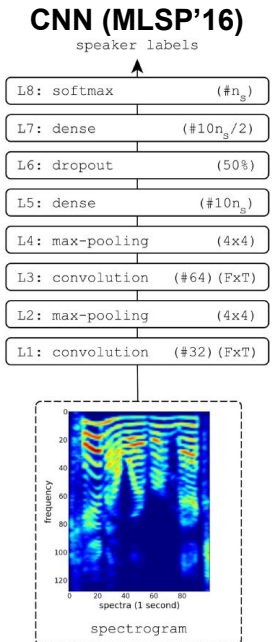
Method	MR	MR (legacy)
RNN /w PKLD	2.19% ($\frac{1.25\%+2.5\%+1.25\%+3.75\%}{4}$)	4.38% (average of 4 runs)
CNN /w PKLD [24]	-	5%
CNN /w cross entropy [23]	-	5%
ν -SVM [40]	6.25%	-
GMM/MFCC [40]	12.5%	-

Lukic, Vogt, Dürr & Stadelmann (2016). «Speaker Identification and Clustering using Convolutional Neural Networks». MLSP'2016.

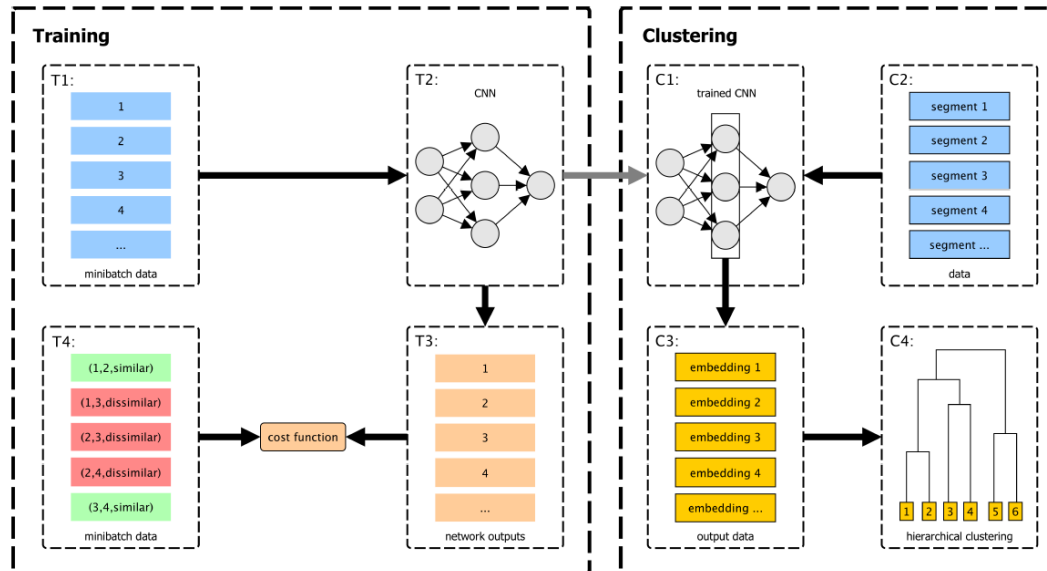
Lukic, Vogt, Dürr & Stadelmann (2017). «Learning Embeddings for Speaker Clustering based on Voice Equality». MLSP'2017.

Stadelmann, Glinski-Haefeli, Gerber & Dürr (2018). «Capturing Suprasegmental Features of a Voice with RNNs for Improved Speaker Clustering». ANNPR'2018.

Speaker clustering – exploiting time information



CNN & clustering-loss (MLSP'17)



Method	MR	MR (legacy)
RNN /w PKLD	2.19% ($\frac{1.25\%+2.5\%+1.25\%+3.75\%}{4}$)	4.38% (average of 4 runs)
CNN /w PKLD [24]	-	5%
CNN /w cross entropy [23]	-	5%
ν -SVM [40]	6.25%	-
GMM/MFCC [40]	12.5%	-

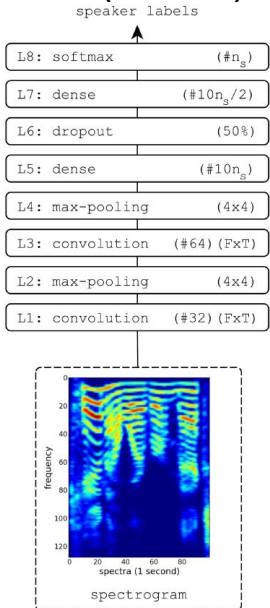
Lukic, Vogt, Dürr & Stadelmann (2016). «Speaker Identification and Clustering using Convolutional Neural Networks». MLSP'2016.

Lukic, Vogt, Dürr & Stadelmann (2017). «Learning Embeddings for Speaker Clustering based on Voice Equality». MLSP'2017.

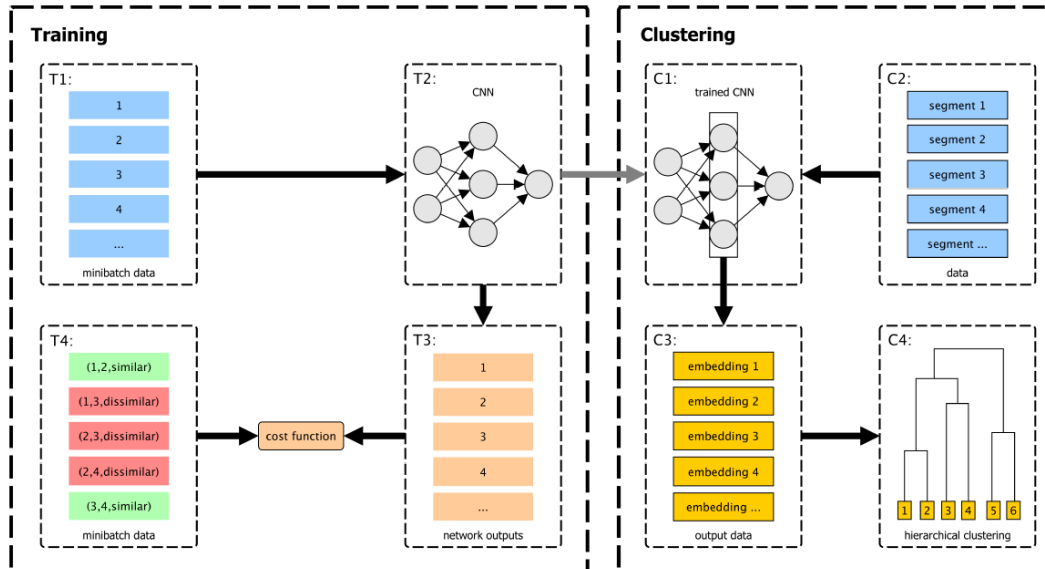
Stadelmann, Glinski-Haefeli, Gerber & Dürr (2018). «Capturing Suprasegmental Features of a Voice with RNNs for Improved Speaker Clustering». ANNPR'2018.

Speaker clustering – exploiting time information

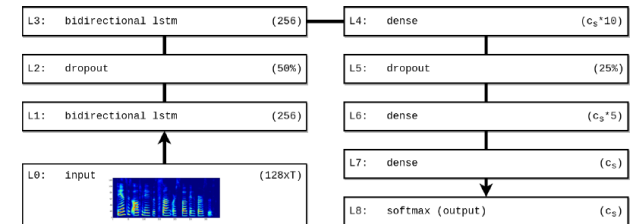
CNN (MLSP'16)



CNN & clustering-loss (MLSP'17)



RNN & clustering-loss (ANNPR'18)



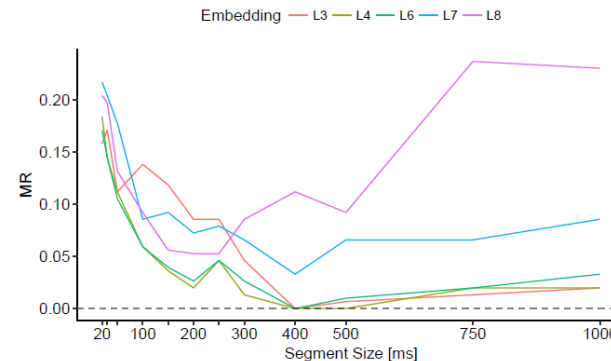
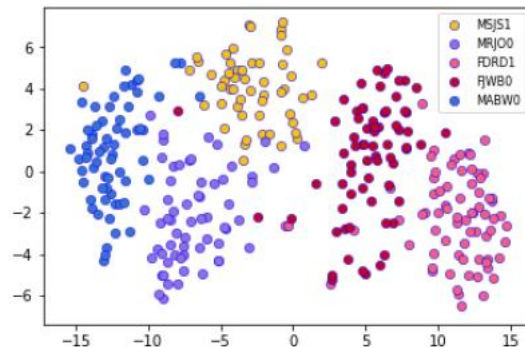
Method	MR	MR (legacy)
RNN /w PKLD	2.19% ($\frac{1.25\%+2.5\%+1.25\%+3.75\%}{4}$)	4.38% (average of 4 runs)
CNN /w PKLD [24]	-	5%
CNN /w cross entropy [23]	-	5%
ν -SVM [40]	6.25%	-
GMM/MFCC [40]	12.5%	-

Lukic, Vogt, Dürr & Stadelmann (2016). «Speaker Identification and Clustering using Convolutional Neural Networks». MLSP'2016.

Lukic, Vogt, Dürr & Stadelmann (2017). «Learning Embeddings for Speaker Clustering based on Voice Equality». MLSP'2017.

Stadelmann, Glinski-Haefeli, Gerber & Dürr (2018). «Capturing Suprasegmental Features of a Voice with RNNs for Improved Speaker Clustering». ANNPR'2018.

Speaker clustering – learnings & future work



«Pure» voice modeling seem largely solved

- RNN **embeddings work well** (see t-SNE plot of single segments)
- RNN model robustly exhibits *the predicted* «**sweet spot**» for the used **time information**
- Speaker clustering on clean & reasonably long input works **an order of magnitude better** (*as predicted*)
- Additionally, using a smarter clustering algorithm on top of embeddings makes **clustering on TIMIT as good as identification** (see ICPR'18 paper on dominant sets)

Future work

- Make models robust on **real-worldish data** (noise and more speakers/segments)
- Exploit findings for robust reliable **speaker diarization**
- **Learn** embeddings and the clustering algorithm **end to end**

Hibraj, Vascon, Stadelmann & Pelillo (2018). «Speaker Clustering Using Dominant Sets». ICPR'2018.

Meier, Elezi, Amirian, Dürr & Stadelmann (2018). «Learning Neural Models for End-to-End Clustering». ANNPR'2018.

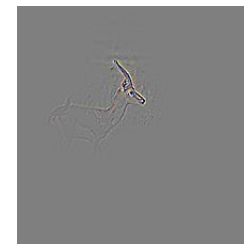
Lessons learned

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)

Robustness is important.

- **Training processes** can be tricky
→ give hints via a unique loss, proper preprocessing and pretraining
- **Risk minimization** instead of error minimization
→ detect all defects at the expense of lower precision



Lessons learned – model interpretability

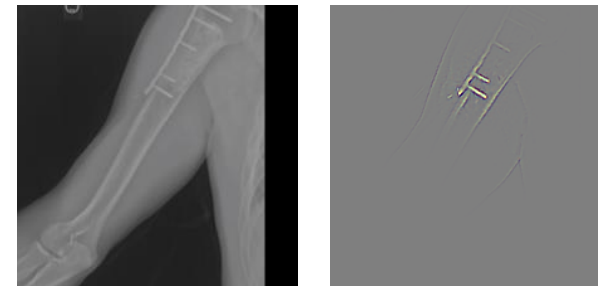
Interpretability is required.

- Helps the developer in «debugging», needed by the user to trust
→ visualizations of learned features, training process, learning curves etc. should be «always on»

negative X-ray



positive X-ray



Stadelmann, Amirian, Arabaci, Arnold, Duivesteyn, Elezi, Geiger, Lörwald, Meier, Rombach & Tuggener (2018). «Deep Learning in the Wild». ANNPR'2018.
Schwartz-Ziv & Tishby (2017). «Opening the Black Box of Deep Neural Networks via Information».
<https://distill.pub/2017/feature-visualization/>, <https://stanfordmlgroup.github.io/competitions/mura/>

Lessons learned – model interpretability

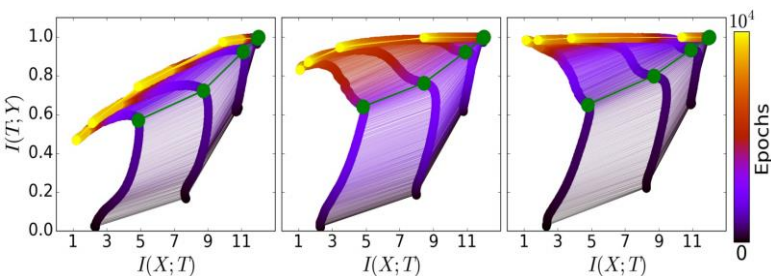
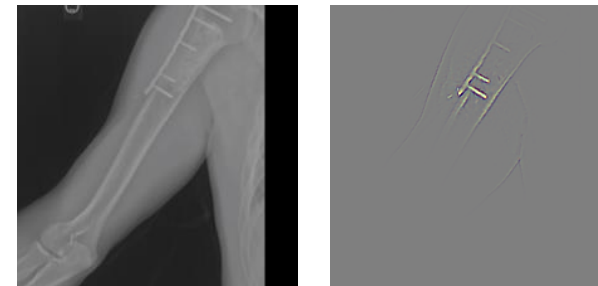
Interpretability is required.

- Helps the developer in «debugging», needed by the user to trust
→ visualizations of learned features, training process, learning curves etc. should be «always on»

negative X-ray



positive X-ray



DNN training on the Information Plane

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

Schwartz-Ziv & Tishby (2017). «Opening the Black Box of Deep Neural Networks via Information».

<https://distill.pub/2017/feature-visualization/>, <https://stanfordmlgroup.github.io/competitions/mura/>

Lessons learned – model interpretability

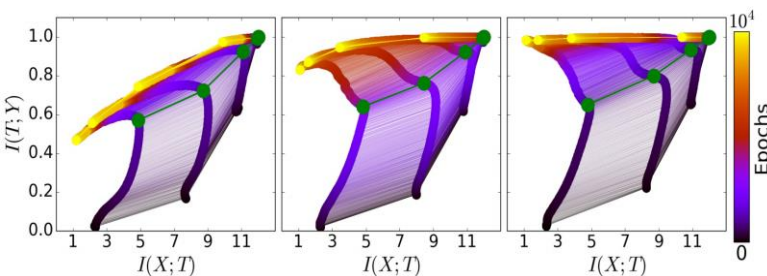
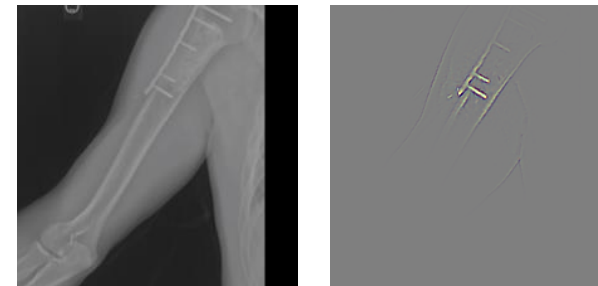
Interpretability is required.

- Helps the developer in «debugging», needed by the user to trust
→ visualizations of learned features, training process, learning curves etc. should be «always on»

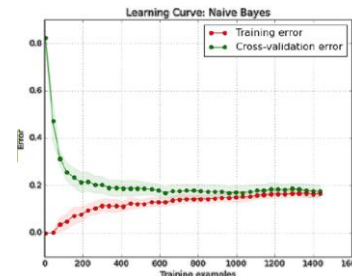
negative X-ray



positive X-ray



DNN training on the Information Plane



a learning curve

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

Schwartz-Ziv & Tishby (2017). «Opening the Black Box of Deep Neural Networks via Information».

<https://distill.pub/2017/feature-visualization/>, <https://stanfordmlgroup.github.io/competitions/mura/>

Lessons learned – model interpretability

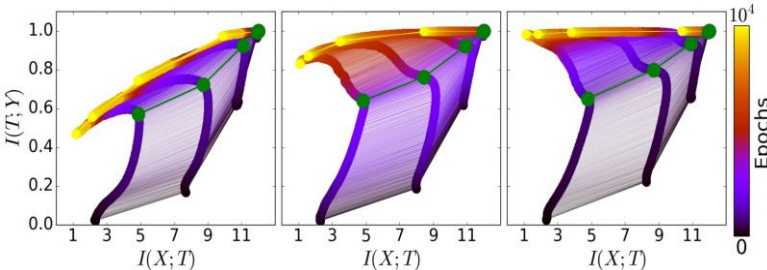
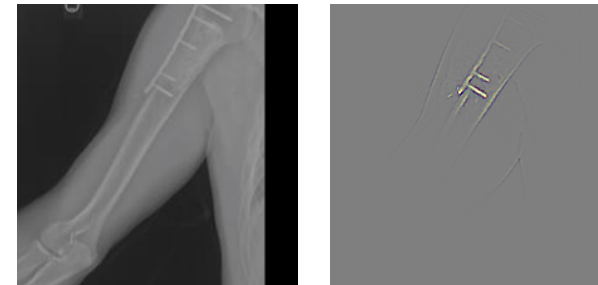
Interpretability is required.

- Helps the developer in «debugging», needed by the user to trust
→ visualizations of learned features, training process, learning curves etc. should be «always on»

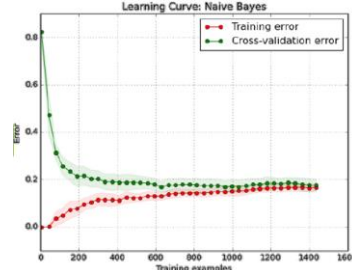
negative X-ray



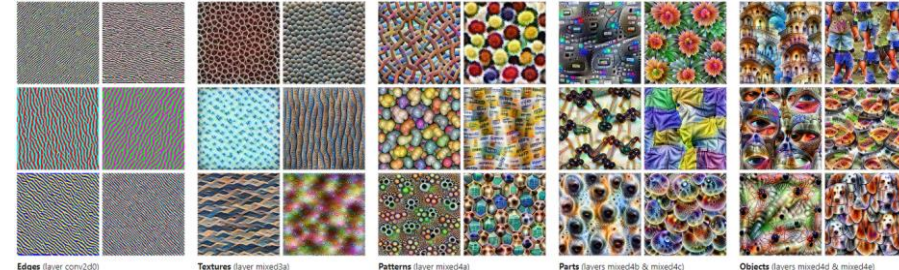
positive X-ray



DNN training on the Information Plane



a learning curve









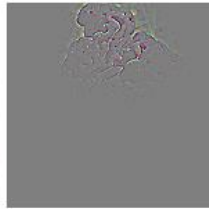
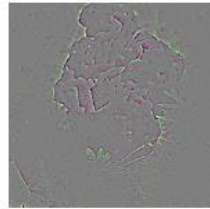
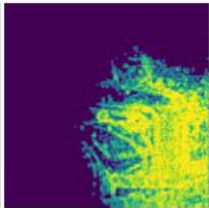
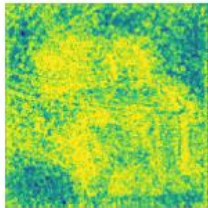
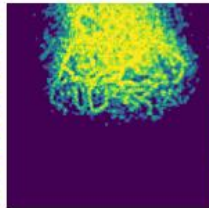
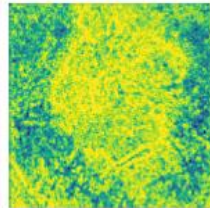
feature visualization

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

Schwartz-Ziv & Tishby (2017). «Opening the Black Box of Deep Neural Networks via Information».

<https://distill.pub/2017/feature-visualization/>, <https://stanfordmlgroup.github.io/competitions/mura/>

Goody – trace & detect adversarial attacks ...using average local spatial entropy of feature response maps

	Original	Adversarial	Original	Adversarial
Image:				
Feature response:				
Local spatial entropy:				

Conclusions

- Deep learning **is applied** and deployed in «normal» businesses (non-AI, SME)
- It does not need big-, but some **data (effort usually underestimated)**
- DL/RL **training** for new use cases **can be tricky** (→ needs thorough experimentation)
- New **theory and visualizations** help to debug & understand
 - *the training process*
 - *individual results*



On me:

- Head ZHAW Datalab, vice president SGAICO, board Data+Service
- thilo.stadelmann@zhaw.ch
- 058 934 72 08
- <https://stdm.github.io/>



On the topics:

- AI: <https://sgaico.swissinformatics.org/>
- Data+Service Alliance: www.data-service-alliance.ch
- Collaboration: datalab@zhaw.ch

→ Happy to answer questions & requests.