Unsupervised Segmentation of Music Symbols using a GAN Detector

Kwon-Young Choi¹ Supervisors : Bertrand Coüasnon¹ Yann Ricquebourg¹ Richard Zanibbi²

¹Intuidoc, Université de Rennes, CNRS, IRISA, F-35000, France ²DPRL, Rochester Institute of Technology, New-York, USA

kwon-young.choi@irisa.fr

2019-06-11





KY Choi-kwon-young.choi@irisa.fr

1 / 28

SegGAN 00000000000000 Conclusion 00



Background

Accidental Detection

SegGAN

4 Conclusion

Optical Music Recognition







Conclusion 00

Optical Music Recognition

- Complex, dense, noisy orchestral scores
- Segmentation problems



Figure: Touching symbols, Broken symbols

Conclusion 00

Optical Music Recognition

- Complex document and notation
- Segmentation problems
 - Importance of good localization!



PhD Focus

Context

- Syntactical recognition for modeling the musical notation and driving deep learning
- Solve segmentation problems with Deep learning for detecting/recognizing symbols



PhD Focus

Objectives

- How to detect/recognize music symbols?
- How to reduce the quantity of annotated data?
 - Syntactical analysis allows generation of subimage with contextual information
 - Bootstrapping samples
 - Synthetic data

Goal

- Build OMR system using a unified set of rule for music notation
- Adapt system to the data using Deep Learning
 - Without requiring heavy manual annotations

Background 00 Accidental Detection

SegGAN 000000000000000 Conclusion

Plan

Background

Accidental Detection

SegGAN

4 Conclusion

Conclusion 00

- Complex music notation \Rightarrow Rule based OMR system
 - Contextual knowledge



Conclusion

- Complex music notation ⇒ Rule based OMR system
 Contextual knowledge
- Segmentation problems \Rightarrow Deep Learning detectors
 - Produce localization + class information



Conclusion

- Complex music notation ⇒ Rule based OMR system
 Contextual knowledge
- Segmentation problems \Rightarrow Deep Learning detectors
 - Produce localization + class information



Conclusion

- Complex music notation ⇒ Rule based OMR system
 Contextual knowledge
- Segmentation problems \Rightarrow Deep Learning detectors
 - Produce localization + class information
- Contextual information: limit the search space



SegGAN 000000000000000 Conclusion

Single Accidental Detection + Rejection



Figure: Dataset examples. In absence of accidental, it's considered as Rejection

Dataset Constitution

Table: Dataset composition

Label	Quantity
No accidental (Reject)	968
Natural	968
Sharp	777
Flat	242
Total	2955

 Square image size: 4x size of interline, resized to 120x120 pixels

- Interline: Vertical distance between two staff lines + height of one staff line
- 2955 samples \Rightarrow Small dataset

Background	

SegGAN 00000000000000 Conclusion 00

Bootstrapping

 \Rightarrow Generate **more data** for localization task.



Figure: Random sampling on different zones

Background	

SegGAN 00000000000000 Conclusion 00

Bootstrapping

\Rightarrow Generate **more data** for localization task.



Figure: Random sampling on different zones

Background	Accidental Detection	SegGAN	Conclusion
00	00●0000	00000000000	00
Bootstrapping			

 \Rightarrow Generate **more data** for localization task. From: 3k samples \Rightarrow to: 25k, 50k, 100k, 200k, 400k data!



SegGAN 000000000000000 Conclusion 00

Intersection over Union

How to evaluate this detection results?



Conclusion 00

Intersection over Union

- How to evaluate this detection results?
 - Compute intersection area



Backg	

Conclusion 00

Intersection over Union

- How to evaluate this detection results?
 - Compute intersection area
 - Compute union area

•
$$IoU = \frac{Inter}{Union}$$



SegGAN 00000000000000

Qualitative Results

0% loU 0%<loU<50% 50%<loU<75% IoU>75%

- Blue is ground-truth
- Red is model prediction

Classic Detection Methods

- Spatial Transformer based detector [Choi, 2019]
- Region based detector like Faster R-CNN [Huang, 2016]
- ...

Background	Accidental Detection	SegGAN	Conclusion
00	000000●	00000000000	00

Results

Table: Results comparing the best Spatial Transformer (ST) based detector, Faster R-CNN, R-FCN and SSD [Huang, 2016]. Results shown are mAP (in %) with an IoU threshold of either 0.5 or 0.75.

Detectors	mAP with IoU > 0.5		mAP with IoU $>$ 0.75	
	μ (%)	$\sigma(\%)$	μ (%)	σ (%)
ST	97.25	1.68	94.81	2.99
Faster R-CNN R-FCN SSD	98.73 99.17 98.93	0.94 0.30 0.67	98.34 98.73 97.81	0.73 0.40 0.92

Conclusion

Plan

Background

Accidental Detection

SegGAN

4 Conclusion

KY Choi-kwon-young.choi@irisa.fr

Problem

Annotation Cost

- Manual annotation for object detection ⇒ costly and time-consuming to produce
- e.g: MUSCIMA++ took 400 hours to produce

Proposal

- Reduce the need for manual data annotation
- Do detection task without explicit detection annotation
- At least generate enough detection annotation in order train a Faster R-CNN

Methodology

Classic Method

- Use synthetic data to train supervised detector
 - Generate synthetic music scores using music typesetting software
 - Train detector (like Faster R-CNN or SSD) using synthetic data
 - Apply trained detector on real not annotated dataset
- Pros:
 - Ability to generate massive amount of data using typesetting software
- Cons:
 - No guaranty that detector trained on synthetic data will work well on real data
 - Synthetic data generation is bound to music typesetting software (maybe there are no equivalent in other domains)

Background	

Methodology

Classic Method

- Use synthetic data to train supervised detector
 - Train detector (like Faster R-CNN or SSD) using synthetic data
 - Apply trained detector on real not annotated dataset
- Cons:
 - No guaranty that detector trained on synthetic data will work well on real data

Proposed Method

- Reduce the input size of images fed to detectors using DMOS grammar (reduce the search space)
- Bridge differences between synthetic data domain and real data domain using Generative Adversarial Network (GAN)

SegGAN 00000000000 Conclusion

Generative Adversarial Network: A Min-Max Game



Figure: Generative Adversarial Network: Discriminator Mode

Discriminator Cost Functions to Maximize

$$\frac{1}{m}\sum_{i=1}^{m}[logD(x^{i}) + log(1 - D(G(z^{(i)})))]$$



SegGAN 00000000000000 Conclusion

Generative Adversarial Network: A Min-Max Game



Figure: Generative Adversarial Network: Adversarial Mode

Adversarial Cost Functions to Minimize

 $\frac{1}{m} \sum_{i=1}^{m} \log(1 - D(G(z^{(i)})))$



Methodology

Key Concept

- Use real isolated music symbol from existing dataset to construct synthetic x training distribution
- Train a Generator to imitate synthetic data distribution
- Detect symbols in images generated by Generator using SSD detector pre-trained on synthetic data

Segmentation GAN

 Replace Generator with a U-Net network [Ronneberger 2015] [Luc 2016]



Figure: SegGan + pre-trained detector

Background	Accidental Detection	SegGAN	Conclusion
		00000000000	
Segmentat	ion CAN		

 Replace Generator with a U-Net network [Ronneberger 2015] [Luc 2016]



Eigure: Using a U. Not for the generate @irisa.fr GAN Detector

 \circ

SegGAN 000000000000

Segmentation GAN

 Replace Generator with a U-Net network [Ronneberger 2015] [Luc 2016]



Figure: Generator should output a mask extracting symbol

Background	Accidental Detection	SegGAN	Conclusion
00		0000●0000000	00
Segmentati	on GAN		

 Replace Generator with a U-Net network [Ronneberger 2015] [Luc 2016]



Figure: Discriminator should discriminate between mask and synthetic image

Background	

Conclusion

Using Isolated Music Symbol Dataset

- Replace ideal x distribution with synthetic data
- Random position and scale using Uniform Distribution
- Arbitrarily set min/max width/height sizes



Figure: Example of synthetic images generated from *real* isolated symbols.

Supervised Generator Loss

Proposition

- Improve training of Generator using additional loss with synthetic data
- Learn identity transform for symbol detection
- Learn rejection task
 - Use non-target class of isolated symbols as rejection examples
 - Define rejection as the task of outputting a blank image

Supervised Generator Loss

 Important: keep the balance in training between Discriminator and Generator



Figure: Supervised Generator Loss

Results

Experimental conditions

RMSProp optimizer

Learning rates balancing Generator and Discriminator training:

- discriminator Ir: 0.0002
- adversarial Ir: 0.0001
- supervised generator loss lr: gridsearch between [1e-5, 1e-3]
- Apply pre-trained SSD detector on generator image
 - Evaluate results using mAP with IoU >0.75 metric

Background	Accidental Detection	SegGAN	Conclusion
00		0000000●0000	00





Background	Accidental Detection	SegGAN	Conclusion
00		0000000●0000	00

Results



Background	Accidental Detection	SegGAN	Conclusion
00		0000000●0000	00

Results



SegGAN 00000000000000

Conclusion



Background 00 Accidental Detectio

SegGAN 00000000000000

Conclusion 00



SegGAN 00000000000000

Conclusion



Backg	

SegGAN 00000000000000

Conclusion



SegGAN 00000000000000

Conclusion



SegGAN 00000000000000

Conclusion



SegGAN 000000000000000

Early Stopping Problem

Instability of GAN training

GAN training is an oscillation around an equilibrium point



KY Choi-kwon-young.choi@irisa.fr

GAN Detector



Early Stopping Methodology

Goal

Find optimal GAN state during training that gives best mAP results

Possible Solutions

- Use synthetic data
 - Current generated synthetic not suitable for early stopping (data is too simplistic)
 - Use synthetic data generated from music typesetting software
 - Possibility to vary iou threshold to better discriminate GAN state
- Use a very small validation set (20/50 examples per class) of manually annotated symbols in real data

Summary

Goal

Train Deep-Learning detectors for music symbol detection using only shape information of an isolated music symbol dataset.

Methodology

- Generate synthetic data using isolated music symbols
- Use SSD detector trained on synthetic data
- Bridge gap between real data and synthetic data using GAN
 - Guide/improve GAN training by designing synthetic data generation strategy and additional loss

SegGAN 00000000000000 Conclusion

Plan

Background

Accidental Detection

SegGAN

4 Conclusion



Conclusion

Overview

- Use GAN model for unsupervised detection of music symbol
 - Bridge differences between synthetic data and real data using Generative model
 - Do detection task without explicit detection annotation (use only shape information from isolated symbols)
 - Reduce the need for manual data annotation

Future Work

- Early stopping of the GAN training
- Improve SegGan model for multi-symbol detection
- Extend application of SegGan to more symbol classes

References

Ronneberger et al. (2015)

U-Net: Convolutional Networks for Biomedical Image Segmentation.

In Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015, 234–41. Lecture Notes in Computer Science

Luc et al. (2016)

Semantic Segmentation Using Adversarial Networks.



Huang et al. (2016)

Speed/Accuracy Trade-Offs for Modern Convolutional Object Detectors



Choi et al. (2019)

CNN-Based Accidental Detection in Dense Printed Piano Scores.

In 2019 15th IAPR International Conference on Document Analysis and Recognition (ICDAR). Sydney, Australia