# Vienna
# Deep Learning
## Meetup

November 20, 2017 @ A1 Telekom Austria

# Vienna
# Deep Learning
## Meetup

## The Organizers:



Thomas Lidy     Jan Schlüter     Alex Schindler

# Vienna
# 14th Deep Learning
# Meetup

**Agenda:**

- Welcome by the Organizers

- Introduction by A1 Telekom Austria (Alexander Stock, Director IT)

- **Evolution of Image Search @ Seznam.cz** (Lukáš Vrabel)

*30 minutes break*

- Hot Topics & Latest News (Tom Lidy, Alex Schindler, Jan Schlüter)

- Discussions and Networking

Vienna
Deep Learning
Meetup

# Hot Topics &
# Latest News

**Tom Lidy,**

**Alexander Schindler,**

**Jan Schlüter**

a short block at every meetup
to briefly present recent papers and news

Send us contributions ([tom.lidy@gmail.com](mailto:tom.lidy@gmail.com))
or come with slides to do a short block yourself!

Vienna
**Deep Learning**
Meetup

# Detection and Classification of Acoustic Scenes and Events

An IEEE AASP Challenge

**DCASE 2017 WORKSHOP**

16 - 17 November 2017, Munich, Germany

**DCASE 2017 CHALLENGE**

15 March 2017 - 31 July 2017

Vienna
**Deep Learning**
Meetup

# Acoustic scene classification

- Classify audio file
- 15 scenes (park, street, office, supermarket, home, etc)
- Multiclass Single Label prediction

Output

Office

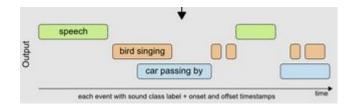# Detection of rare sound events

- rare sound events in artificially created mixtures
- 3 classes (Baby Crying, Glass Breaking, Gunshot)
- Sequence labelling (start/stop of event), one class per file

Output

glass breaking

each event with sound class label + onset and offset timestamps    time

# Sound event detection in real life audio

- Classify audio file
- 6 classes (brakes squeaking, car, children, people speaking, people walking)
- Sequence labelling (start/stop of event), multi-label classification

Output

speech

bird singing

car passing by

each event with sound class label + onset and offset timestamps    time

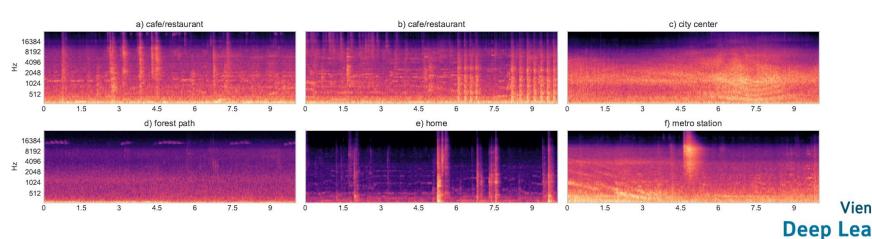# Large-scale weakly supervised sound event detection for smart cars

- subset of "AudioSet: An Ontology And Human-Labeled Dataset For Audio Events"
- 17 sound events divided into two categories: "Warning" and "Vehicle"
- Sequence labelling (start/stop of event), multi-label classification

Vienna
**Deep Learning**
Meetup

# Multi-resolution Convolutional Neural Networks

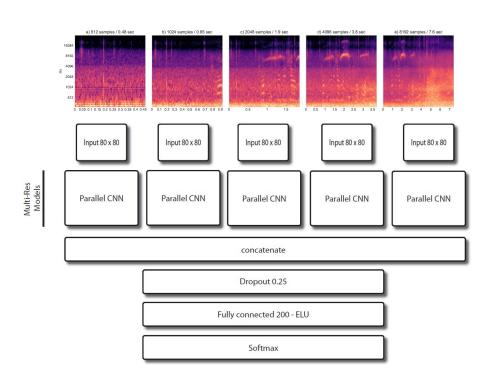Alexander Schindler, Thomas Lidy and Andreas Rauber

- **Problem:** Observation DCASE 2016
  - Confusion between Train / Tram / Bus / Supermarket
  - Low-Frequent Humming
    - Diesel engine, air-condition, refrigerators
  - Acoustic scene composed of
    - Timbral texture (short-term)
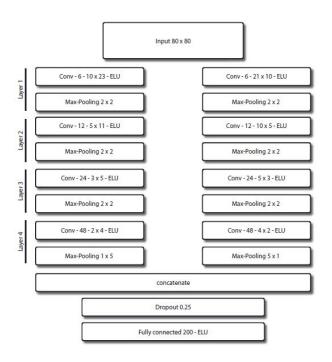    - Acoustic events (long-term)

# Multi-resolution Convolutional Neural Networks

Alexander Schindler, Thomas Lidy and Andreas Rauber

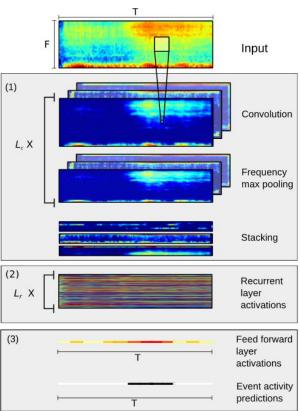- **Solution:** Train Model on multiple-temporal resolutions

# DCASE Workshop
## Hot Topics in
## Audio event Detection

- **Recurrent Convolutional Neural Networks** (RCNN)
  - dominant in audio event detection
  - CNN learns audio representation
  - CNN processes sequential slices of Spectrogram input
  - Stacking CNN outputs
  - Recurrent Layer
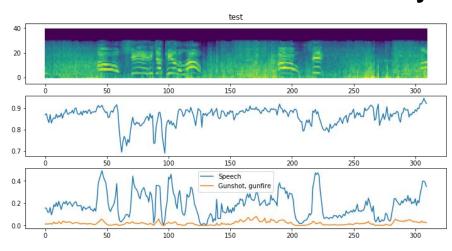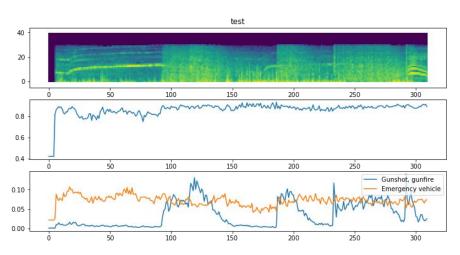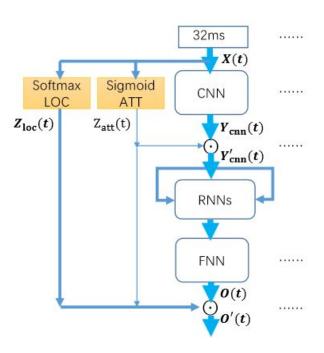  - Sequential ("TimeDistributed") Fully Connected Layers

# DCASE Workshop
## Hot Topics in
## Audio event Detection

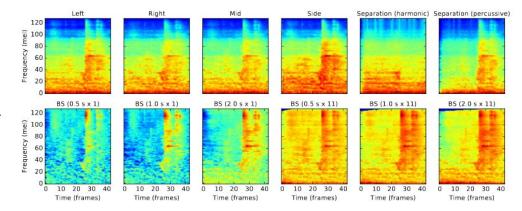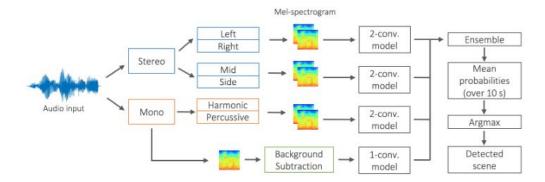- **Attention & Localization Layers**

# DCASE Workshop
## Hot Topics in
## Audio event Detection

- **Binaural (stereo) Models**
  - 2nd best model
  - Uses all input channels
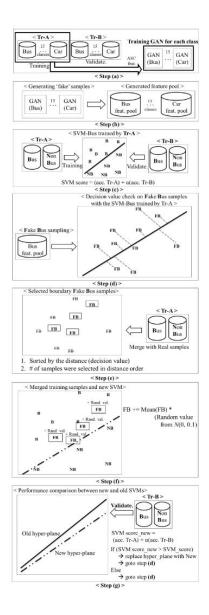    - (usually mono input)

# DCASE Workshop
## Hot Topics in
## Audio event Detection

- **Generative Adversarial Network (GAN) based Audio Data Augmentation**
  - <u>Winning Model Task 1</u>
  - Uses Gan to generate training examples
  - Generates feature vectors - not audio

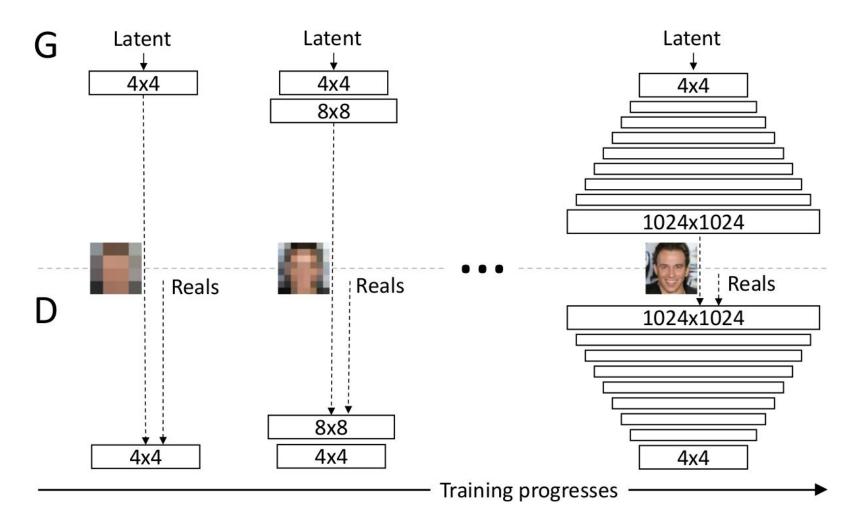# Deep Voice 3: 2000-Speaker Neural Text-to-Speech

- Trained on LibriSpeech dataset (820 hours, 2484 speakers)
- Similar quality to Tacotron, but faster
- No comparison to DeepMind's new WaveNet (which has no paper yet)
- Interesting evaluation:

| Text Input | Attention | Inference constraint | Repeat | Mispronounce | Skip |
|---|---|---|---|---|---|
| Characters-only | Dot-Product | Yes | 3 | 35 | 19 |
| Phonemes & Characters | Dot-Product | No | 12 | 10 | 15 |
| **Phonemes&Characters** | **Dot-Product** | **Yes** | **1** | **4** | **3** |
| Phonemes & Characters | Monotonic | No | 5 | 9 | 11 |

Number of errors made in vocalizing a set of 100 test sentences

http://arxiv.org/abs/1710.07654

Vienna
**Deep Learning**
Meetup

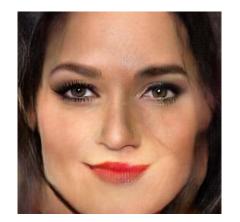# Progressive Growing of GANs for Improved Quality, Stability, and Variation

# Progressive Growing of GANs for Improved Quality, Stability, and Variation

Existing work:

This work:

Vienna
**Deep Learning**
Meetup

# Progressive Growing of GANs for Improved Quality, Stability, and Variation

But:





https://github.com/tkarras/progressive_growing_of_gans
http://arxiv.org/abs/1710.10196

Vienna
**Deep Learning**
Meetup

# Optimal transport maps for distribution preserving operations on latent spaces of Generative Models

- GAN learns mapping from random vector in latent space to image
- Typical demonstration in paper: Interpolate in latent space
- **Problem:** in high dimensions, volume of a sphere/cube concentrates near surface/edges ⇒ interpolation traverses space hardly seen in training
- This can be fixed!



(b) Uniform prior distribution.



(c) Linear midpoint distribution



(d) Matched midpoint distribution (**ours**)



(e) Spherical midpoint distribution (White, 2016)

https://arxiv.org/abs/1711.01970

Vienna
**Deep Learning**
Meetup

# mixup: Beyond Empirical Risk Minimization

- Proposed regularization method:
  - Draw factor λ from beta distribution (α = β = 0.2; so mostly near 0/1)
  - Input: linear mixture of two data points $x = \lambda x_1 + (1-\lambda)x_2$
  - Target: linear mixture of respective targets $y = \lambda y_1 + (1-\lambda)y_2$
    (with $y_1$ and $y_2$ represented as one-hot vectors)

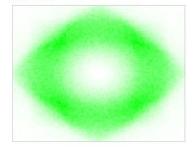- Improved performance on CIFAR-10/100, ImageNet, Google Commands

| Dataset | Model | ERM | *mixup* |
|---|---|---|---|
| CIFAR-10 | PreAct ResNet-18 | 5.6 | **3.9** |
| | WideResNet-28-10 | 3.8 | **2.7** |
| | DenseNet-BC-190 | 3.7 | **2.7** |
| CIFAR-100 | PreAct ResNet-18 | 25.6 | **21.1** |
| | WideResNet-28-10 | 19.4 | **17.5** |
| | DenseNet-BC-190 | 19.0 | **16.8** |

(a) Test errors for the CIFAR experiments.



(b) Test error evolution for the best ERM and *mixup* models.

https://arxiv.org/abs/1710.09412
http://www.inference.vc/mixup-data-dependent-data-augmentation/

Vienna
**Deep Learning**
Meetup

# ImageNet in 15 minutes

- ResNet-50 on ImageNet
- Batchsize: 32k
- Hardware: 1024 NVIDIA P100 (cost: EUR 6000 each)
- Training time: 15 min for 90 epochs

| Team | Hardware | Software | Minibatch size | Time | Accuracy |
|------|----------|----------|----------------|------|----------|
| He *et al.* [5] | Tesla P100 × 8 | Caffe | 256 | 29 hr | 75.3 % |
| Goyal *et al.* [4] | Tesla P100 × 256 | Caffe2 | 8,192 | 1 hr | 76.3 % |
| Codreanu *et al.* [3] | KNL 7250 × 720 | Intel Caffe | 11,520 | 62 min | 75.0 % |
| You *et al.* [10] | Xeon 8160 × 1600 | Intel Caffe | 16,000 | 31 min | 75.3 % |
| This work | Tesla P100 × 1024 | Chainer | 32,768 | 15 min | 74.9 % |

- Some tricks (mostly from Goyal et al.):
  - warm-up with low learning rate
  - warm-up with RMSprop, then switch to SGD
  - no moving averages in Batch Normalization

https://www.preferred-networks.jp/en/news/pr20171110

Vienna
**Deep Learning**
Meetup

# AI Robot Learned How to Pick up Objects

- University of California, Berkeley, trained a deep learning system on a **simulated** dataset of more than a thousand objects
- DNN exposed to each one's 3D shape and appearance
- Test on <u>physical</u> objects that weren't included in its <u>digital</u> training set

- When the system thought it had a better than 50 percent chance of successfully picking up a new object, it was actually able to do it **98 percent** of the time
  — without having trained on any real-world objects
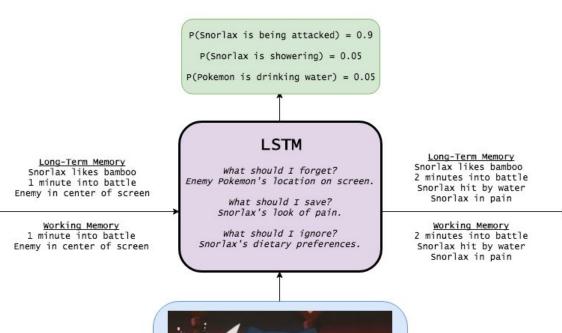- shows that a simulated data set can be used to train a model for grasping

Vienna
**Deep Learning**
Meetup

# AI Robot Learned How to Pick up Objects

Vienna
**Deep Learning**
Meetup

# Exploring LSTMs



P(Snorlax is being attacked) = 0.9

P(Snorlax is showering) = 0.05

P(Pokemon is drinking water) = 0.05

## LSTM

what should I forget?
Enemy Pokemon's location on screen.

what should I save?
Snorlax's look of pain.

what should I ignore?
Snorlax's dietary preferences.

**Long-Term Memory**
Snorlax likes bamboo
1 minute into battle
Enemy in center of screen

**Working Memory**
1 minute into battle
Enemy in center of screen

**Long-Term Memory**
Snorlax likes bamboo
2 minutes into battle
Snorlax hit by water
Snorlax in pain

**Working Memory**
2 minutes into battle
Snorlax hit by water
Snorlax in pain
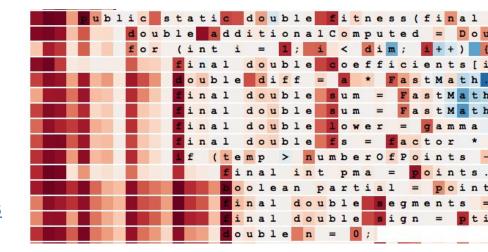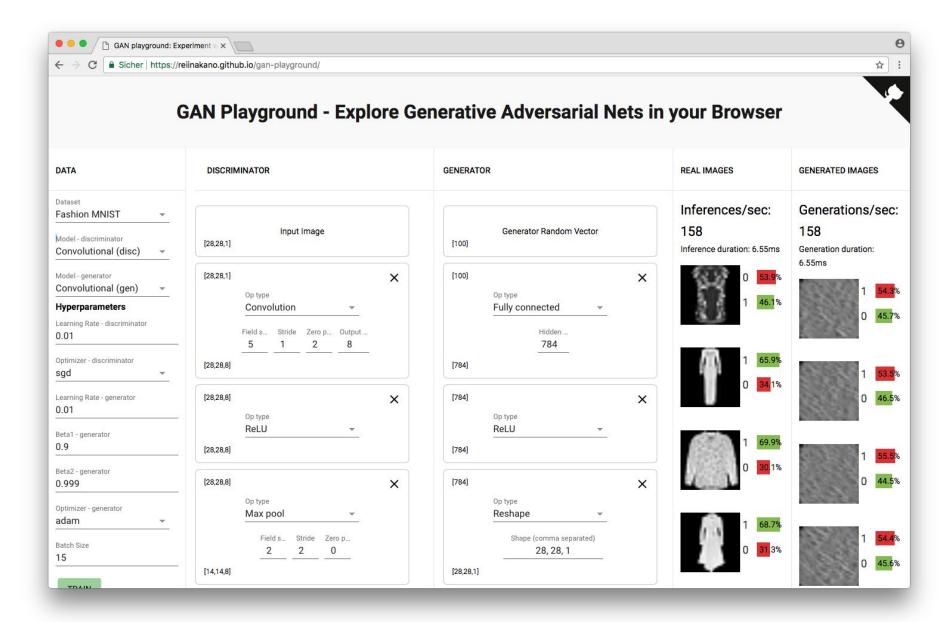
http://blog.echen.me/2017/05/30/exploring-lstms

https://reiinakano.github.io/gan-playground/
(Chrome only)

# Announcements

# Post-doc Position



**Who are we?**
Pharmacoinformatics research Group of Prof. Ecker

**What are we doing?**
Machine learning/ Deep learning for toxicity of chemicals/drugs

**We need you!**
- PhD (preferably in life sciences)
- Experience in Machine learning/ Deep learning
- Scripting/Programming experience: at least 1 language (Python, R,...)

**Page: https://pharminfo.univie.ac.at/**
**Mailto: gerhard.f.ecker@univie.ac.at**

# Ethics and Bias in DL / AI

- Planned for DL Meetup February or March
- Please send us contributions:
  - Papers
  - Articles
  - Speakers
  - Panelists
  - Questions?



Vienna
**Deep Learning**
Meetup

# One more Hot Topic...

New Deep Learning PhD-Thesis !

Vienna
**Deep Learning**
Meetup

# Deep Learning for Event Detection, Sequence Labelling and Similarity Estimation in Music Signals

Doctoral Thesis
to obtain the academic degree of
Doktor der technischen Wissenschaften
in the Doctoral Program
Technische Wissenschaften

- **60-page introduction** to **deep learning** for audio processing
- extended versions of **publications achieving state-of-the-art result** in
  - music/speech detection
  - onset detection
  - segmentation
  - singing voice detection
  - acceleration of music similarity estimation
- previously unpublished negative results from follow-up experiments that didn't work out for each of these tasks
- source code for the singing voice detection experiments
  - https://github.com/f0k/ismir2015
  - https://github.com/f0k/ismir2015/tree/phd_extra
- source code for the wiggly lines drawing demo of ISMIR 2016
  - https://github.com/f0k/singing_horse

Vienna
**Deep Learning**
Meetup

Next

Vienna
**Deep Learning**
Meetup

9th January 2018

weXelerate

Thomas Lidy      Jan Schlüter      Alex Schindler