# CNN-Based Malware Visualization and Explainability

**Lara Dedic**
**Applied Machine Learning Researcher**
**ldedic@novetta.com**
**www.novetta.com**

# Malware Analysis

**Challenges**

- Complex malware is inspected manually
- Takes a long time and lots of effort for a single executable to be analyzed
- Extensive domain knowledge is required

# Machine Learning

**Malware Classification**

- "Deep Learning on Disassembly Data." Davis & Wolff. Black Hat 2015
- "Activation Analysis of a Byte-based Deep Neural Network for Malware Classification" Coull. CAMLIS 2018
- "Malware Detection by Eating a Whole EXE" Raff, Barker, Sylvester, Brandon, Catanzaro & Nicholas. arXiv 2017

**Explainable ML**

- Activation Atlases - Google & OpenAI
- InterpretML - Microsoft

# Proof Of Concept

**Input**
Executable → Image

**CNN**
ResNet

**Output**
Grad-CAM Heatmaps



Malware
Classification

# Convolutional Neural Network

**fast.ai**

- Simplifies applying state-of-the-art deep learning models and techniques
  - Enables rapid prototyping
- Deep Learning Python Library & MOOC
  - Top Down Approach
  - PyTorch base

```
cnn_learner(data, models.resnet50, metrics=[accuracy])
```

# Grad-CAM

**Gradient-weighted Class Activation Mappings**

- Generated from class-specific gradients passed to a convolutional layer of a CNN
  - Retain spatial information
  - Deeper layers capture higher level visual constructs
- No retraining required
- Guided Grad-CAM: Class discriminative and high resolution



Grad-CAM for "Cat"    Grad-CAM for "Dog"

# Grad-CAM as Visual Explanations

**Establish trust in our models**

- Debug ML
- Detect Bias
- Explain unexpected predictions

**Machine teaching**

- Let models teach humans how to make better decisions about data

**Leverage it in other domains**

- Medicine
- Malware reverse engineering and analysis

# Setup

## Data

- ~1400 Windows PEs
- 70% benignware , 30% malware

## Training

- Train ResNet-50 for malware vs benignware classification
- Grad-CAM to be generated from final Conv Layer
- AWS p2.xlarge

## Results

- F1 Score: 96.4%

Confusion matrix

|  | benignware | malware |
|---|---|---|
| benignware | 95 | 3 |
| malware | 4 | 39 |

Actual / Predicted

# APT1 Heatmaps



BOUNCER  AURIGA  TARSIP-ECLIPSE  GREENCAT  GOGGLES  WEBC2-UGX  GLOOXMAIL

# Challenges

**Image Representations of Malware**

- Variable length images
- Code execution not represented well in this format

**Evaluating Heat Maps**

- Human judgement

**Small Dataset**

- ~ 1400 samples from VirusTotal

# SpecAugment
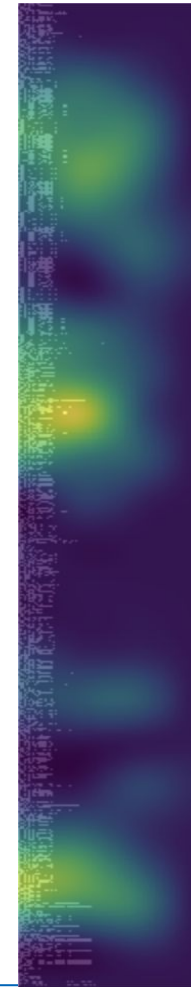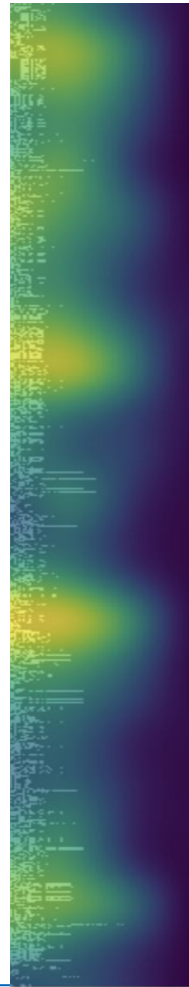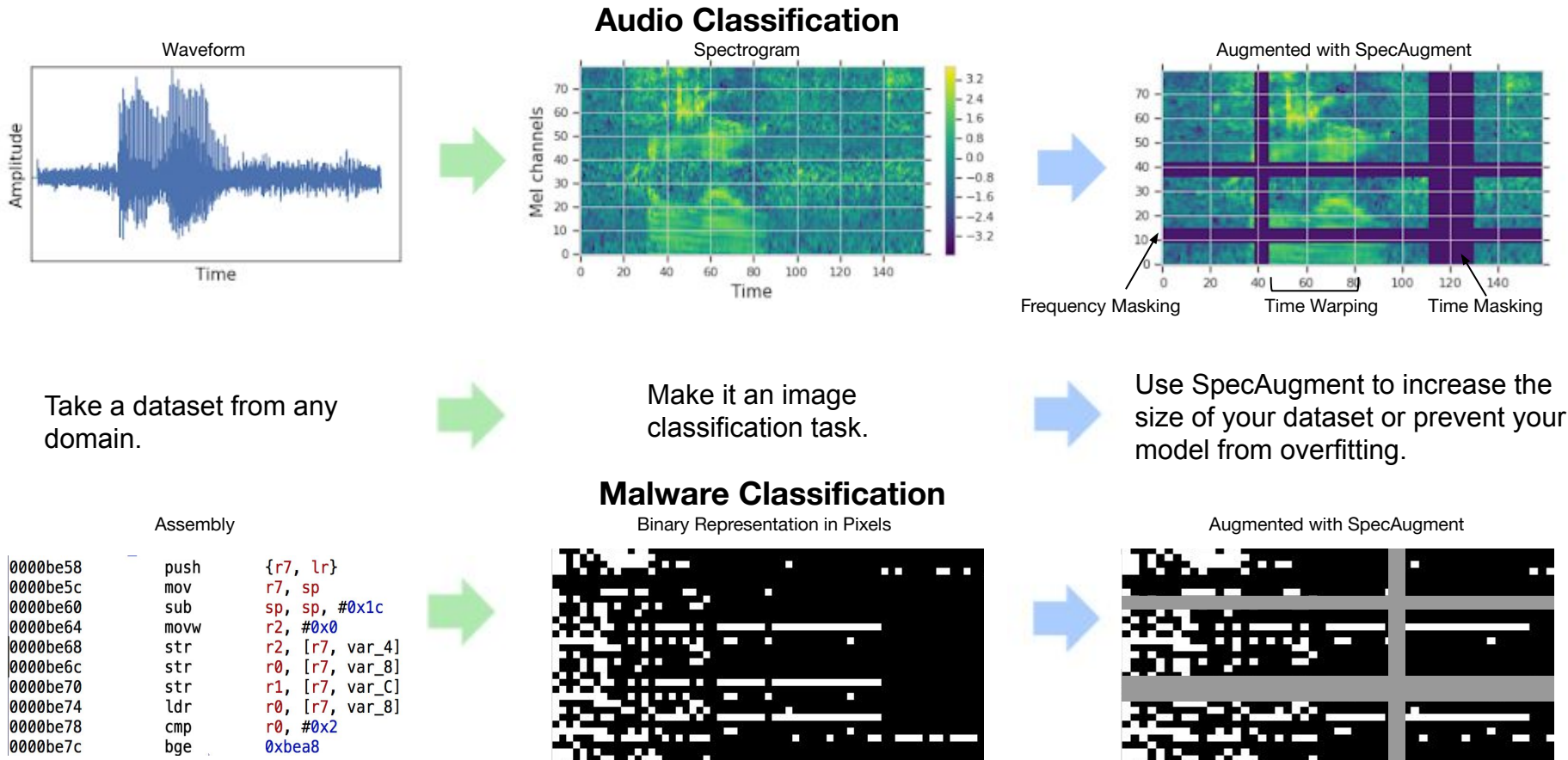
SpecAugment is a state of the art data augmentation technique created by Google Brain in April of 2019 for automatic speech recognition tasks.

## Audio Classification



Waveform

Spectrogram

Augmented with SpecAugment

Frequency Masking     Time Warping     Time Masking

Take a dataset from any domain.

Make it an image classification task.

Use SpecAugment to increase the size of your dataset or prevent your model from overfitting.

## Malware Classification

Assembly

```
0000be58    push    {r7, lr}
0000be5c    mov     r7, sp
0000be60    sub     sp, sp, #0x1c
0000be64    movw    r2, #0x0
0000be68    str     r2, [r7, var_4]
0000be6c    str     r0, [r7, var_8]
0000be70    str     r1, [r7, var_C]
0000be74    ldr     r0, [r7, var_8]
0000be78    cmp     r0, #0x2
0000be7c    bge     0xbea8
```

Binary Representation in Pixels

Augmented with SpecAugment

# What's Next

- **Malware Family Classification**
  - Use Guided Grad-CAM to help analysts discover characteristics of classes of malware
- **Larger Dataset**
- **Input Configuration**



Guided Grad-CAM for "Cat"

Guided Grad-CAM for "Dog"

# Links

- **[Grad-CAM](#)**
- **[fast.ai](#)**
- **[SpecAugment](#)**
- **[InterpretML](#)**
- **[Activation Atlases](#)**