# ml-audio-dev-tools

## Development tools for deep learning models of acoustical signal processing
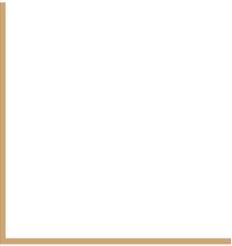
### Scott H. Hawley, Ph.D.

@drscotthawley

- ❖ Chem. & Phys. Dept, Belmont U.
- ❖ Belmont Data Collaborative
- ❖ Harmonai

Brought to you by ASA Technical Committee on Signal Processing
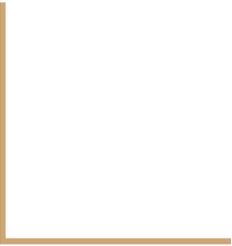
# ml audio dev tools

...are *quite* few in number
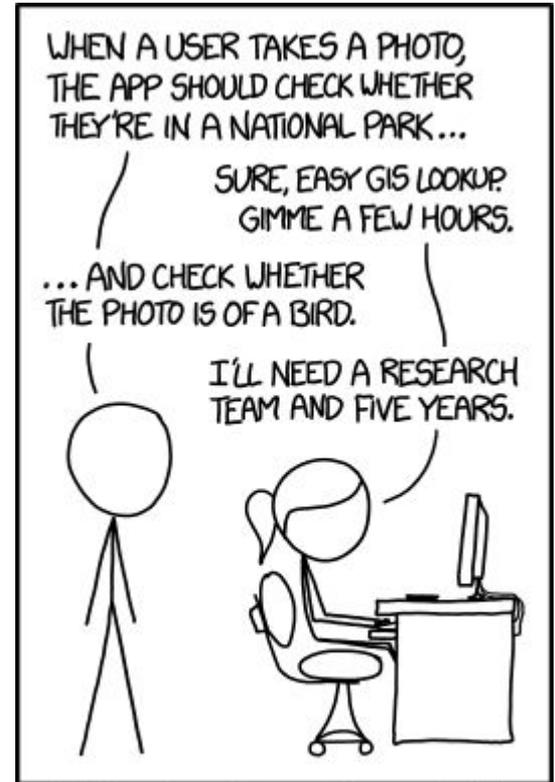
...so we should all build more.

# What's too difficult?

[Jeremy Howard](#) of [fast.ai](#) uses this [XKCD](#) cartoon at the start of Lesson 1 in the fast.ai course "[Practical Deep Learning for Coders](#)"

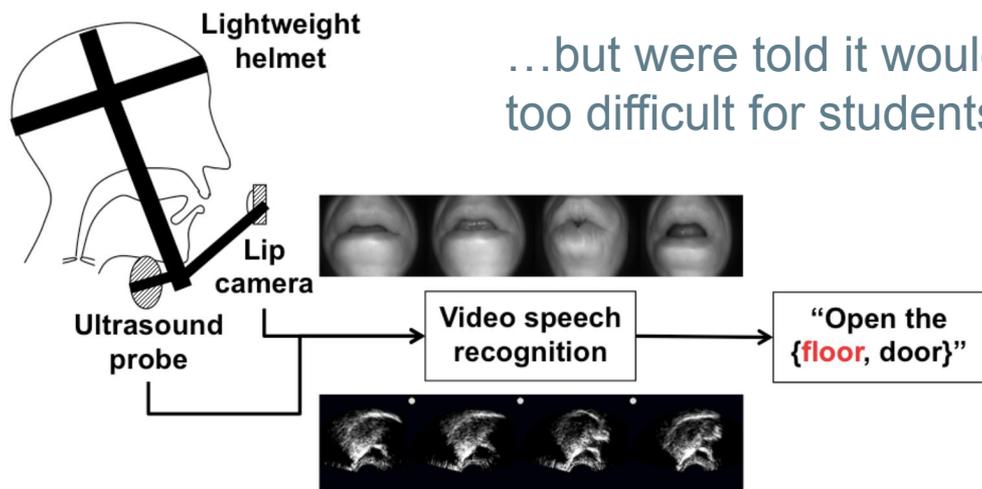...to illustrate the **huge change in the capabilities of Computer Vision systems since the advent of Deep Learning**.

This problem of bird detection is *trivial* nowadays. It is not even worth a homework problem.

# What's too difficult?

In 2020, Bruce Denby (Inst. Lagevin, Sorbonne) & I wrote to CA TC with a proposal for a 2021 ASA Student Challenge in Speech-To-Text using video imagery ("Silent Speech"):
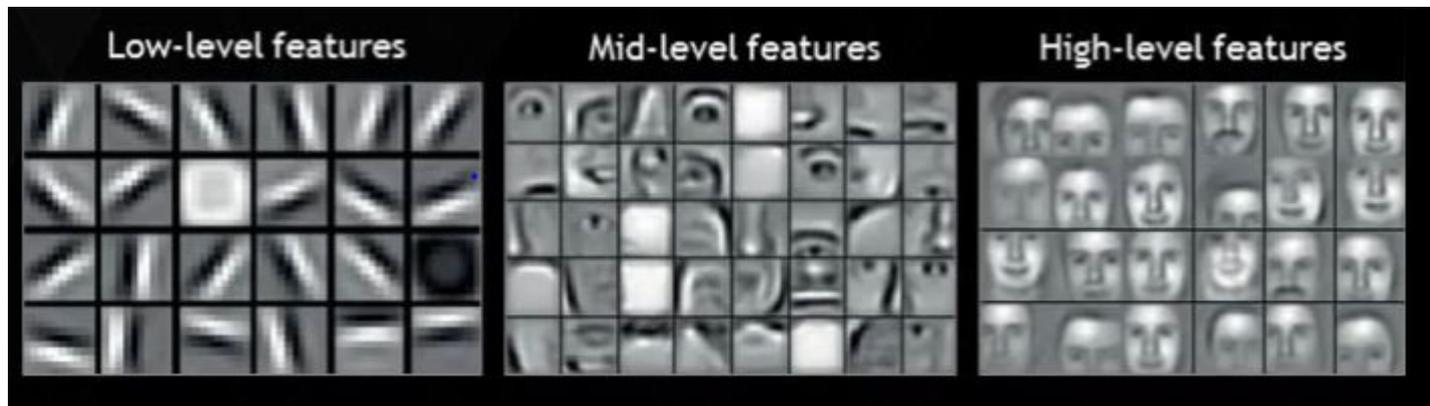
…but were told it would be too difficult for students.

We were *confused.*
With modern DL-CV tools,
This is a *homework problem,*
or at best a minor Kaggle competition.

# *Some* things that were hard become easy with DL

Many human-designed **feature detection** systems can be obviated with Deep Learning systems, which build their own feature detectors as part of the training process.



*An example hierarchy of feature detectors for a facial recognition system. (Source: Lee et al, 2009.)*

*Aside: At SP TC meeting Tuesday night, there was talk about a future session on "Feature Extraction and Dimensionality Reduction" – both of these are built in to DL systems.*

# *Some* things that were hard become easy with DL

Various examples, but here's a couple:

Neural networks are great at *pattern matching* and *denoising*:

I wrote my first denoising autoencoder before I knew what a Wiener filter was!

I wrote an object detector for ellipses (Hawley & Morrison JASA 2021 & JASA Express Letters 2022) without bothering to try Hough Transforms (b/c it was complex)

# Which means...

- there's a new generation of coder-scientists (🙋‍♂️) who **hope** that DL will make up for their *lack of signal processing domain knowledge*. **Not without reason:** DL has shown to *beat* former baselines in many fields for many problems.
- there are veterans in the SP field who are perhaps not up to speed on the rapid pace of advancements in DL for "audio", and may be curious about incorporating DL into their work.  ...Good news: This is/(can be) "easy"*!      *compared to hard-core SP/math, if you're already good at coding

More good news: "Everything old is new again". Much of classic SP still finds its way into DL systems & helps drive innovation. (e.g., Vector Quantization!)

# Who this talk is for:

- Students and ❤️*Esteemed SP Experts*🙇‍♂️: "Onboarding"
- Experienced DL-Audio researchers: Sharing tips!

# Confession/bias of mine:

- When I say "audio" I usually mean "musical" audio:
  - multi-channel
  - high sample rates (44.1kHz+)
  - ...and *not just classification problems*
  - I tend to not even think about *"Speech". No offense*
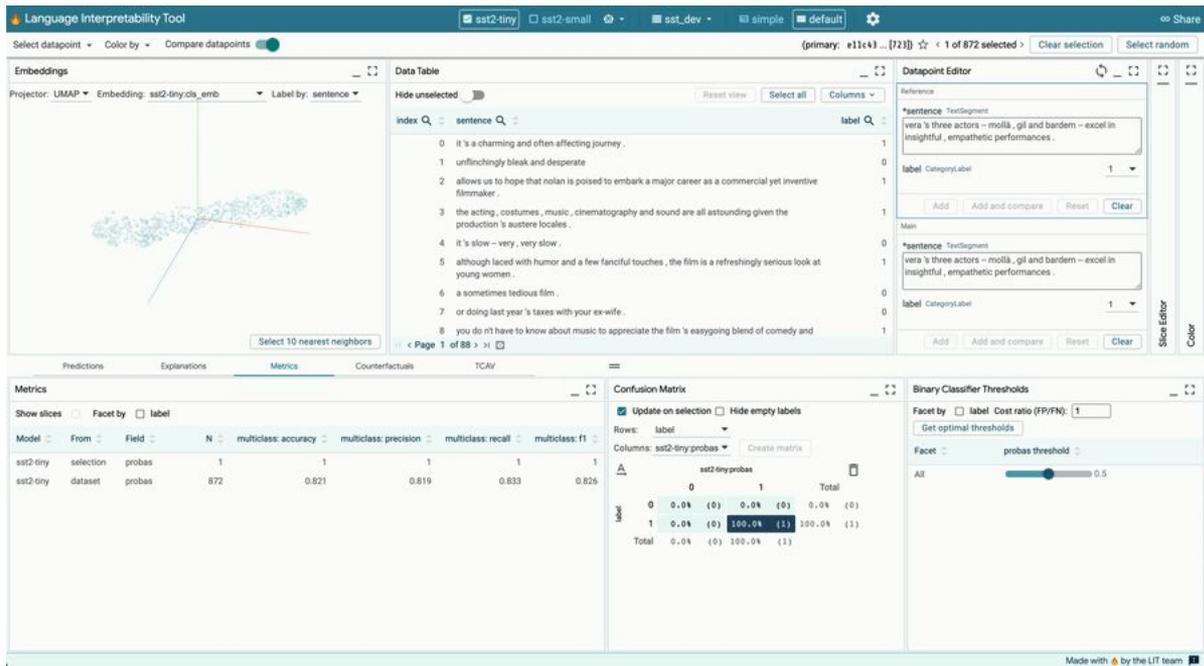- For me, DL ⊂ ML

ml audio dev tools...

...are few in number

# ...**Compared** to Tools for DL Models for **Images**, **Text**, ...and even *Speech*

In these other domains, DL is very mature.*

MANY, many, many tools & demos for viz., analysis, saliency,...

*because money



There's no audio tool as cool as Google PAIR's Language Interpretability Tool

# Responses I got...

i.e., tools suggested by DL-audio practitioners...
Thanks: Brian McFee, Fabian-Robert Stöter, Jesse Engel, Gene Kogan, Dadabots, Andrew Parker, David Braun, Zach Evans, scart97, Eric Hallahan, & Christian Steinmetz

...were mostly:

I. Basic ML workflow / 101
II. Normal audio workflow
III. A few extra cool tools



Scott H. Hawley @drscotthawley · 4d
ML-audio people! What are your favorite "tools" (interpreted very broadly, however you like) for doing development/research work in ML-audio? I want to compile a list for an "onboarding" talk late next week @acousticsorg. Will gladly attribute you for your suggestion & even...

💬 2    🔁 2    ♡ 4

Scott H. Hawley @drscotthawley · 4d
(time permitting) showcase your work. "Tools" could be for coding, testing, viz/auralization, communication, libraries, IDE-fu, computational tricks,...anything you think is handy that others should know about. Audience will include bio & underwater acous. ppl, not just music.

# I. Basic ML - GPU Computing

Graphics Processing Units (GPUs)
are key: 100x+ faster than CPU

You don't need to buy a GPU.
Various cloud-based systems let you
do GPU computing for (near) free.

**Google Colab**, Kaggle, Paperspace
Gradient, Amazon Sagemaker,...

These typically make use of Jupyter notebooks.

(For awesome lib-dev via Jupyter notebooks, check out nbdev.)

# I. Basic ML - Python Programming

Python is overwhelmingly the most popular & well-supported  programming language for ML. (Much more so than MATLAB, C++, Julia, R, JS,...)

There's usually a library/package that does what you want.

Which DL Library? ➡️**PyTorch**⬅️, Tensorflow, JAX,..
               **Lightning**

I hate `conda` and only use venv+ `pip`

# I. Python Packages for Audio and/or DL

- [librosa](): "Swiss army knife":   (Also by Brian McFee: [mir_eval]())
- [torchaudio](): GPU processing
- [auraloss](): Loss functions                    (Christian Steinmetz)
- [Pedalboard](): Data augmentation    (Peter Sobot/Spotify)
- [DawDreamer](): Python DAW          (David Braun)
- [ONNX](): Export to JUCE / C++ / JS

# II. Normal Audio Workflow

- Audacity ← We'll come back to this one!
- Sonic Visualizer:
- Reaper
- Logic



Suggested by Jesse Engel

# II. Audio workflow – VSCode [audio-preview](#)

(Suggested by
Fabian-Robert Stöter)



Works in Remote mode!
i.e., play on your laptop
the files on your server.

# II.5 Trick: (Mel-)Spectrograms + Image-based DL

- Many posts on "The remarkable effectiveness of Convolutional Neural Networks on (Mel-)Spectrograms"
- Translation equivariance of CNNs fits well with phase (+ pitch) translation invariance of human auditory system.
- Upshot: Just using images of spectrograms with an image-based code can work surprisingly well.
- Makes a great baseline before going for full end-to-end audio DL



Source: LibRosa (McFee et al)

If you're just starting out, probably start with this method

# III. Extra Cool – WandB (Audio) Callbacks!

- Weights & Biases ("wandb") is a cloud-based data-logging service you can use for free.

  (Hawley & Morrison JASA-EL 2022 found it "essential" for keeping track of many, many runs.)

- Among the things you can log & playback are audio examples:



Tables Tutorial: Recreating Whale Melodies on Orchestral Instruments

Interactively exploring ML data and predictions in the audio domain with our new Tables feature

| | Whale song | Audio features | Spectrogram | Species | Location |
|---|---|---|---|---|---|
| 1 | | | | Bowhead Whale | Barrow, Alaska CC2A |
| 2 | | | | Bowhead Whale | Bailie Is., Beaufort Sea X |
| 3 | | | | Humpback Whale | St. David's Island, Bermuda |

# III. Communication: Demo Hosting via Gradio.app

Demo for Christian Steinmetz & Josh Reiss, NeurIPS 2021

# III. Label Studio

# III. Extra Cool – Audacity DL Models!

- Download models off 🤗

  (& upload your own first!)

- Run model as audio effect

- So far only mono, no knobs

Note: This is a custom Audacity build, download from [here](#)

*Help them add sliders!*



INTERACTIVE AUDIO LAB    Projects    Publications    Resources    People
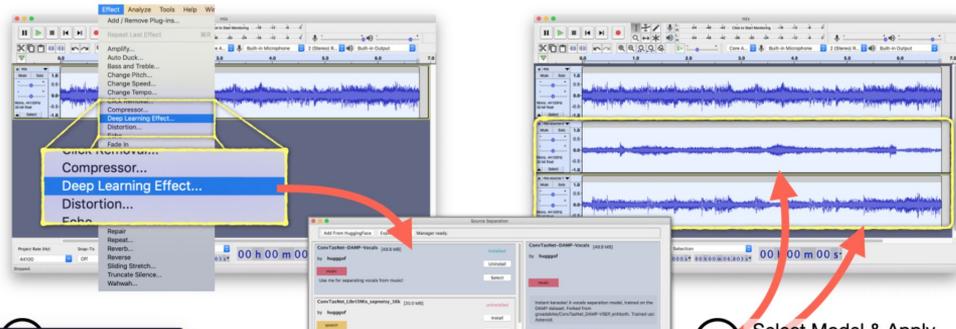
## Deep Learning Tools for Audacity

Hugo Flores Garcia, Aldo Aguilar, Ethan Manilow, Dmitry Vedenko and Bryan Pardo

We provide a software framework that lets deep learning practitioners easily integrate their own PyTorch models into the open-source Audacity DAW. This lets ML audio researchers put tools in the hands of sound artists without doing DAW-specific development work.

Our software framework lets ML developers easily integrate new deep-models into Audacity, a free and open-source DAW that has logged over 100 million downloads since 2015. Developers upload their trained PyTorch model to HuggingFace's Model Hub. The model becomes accessible through Audacity's UI and loads in a manner similar to traditional plugins.

# III. MPF tools like Essentia(.js)

| Similarity | Classification | Deep learning inference | Mood detection |
|---|---|---|---|
| Analyze audio and compute features to find similar sounds or music tracks. | Classify sounds or music based on computed audio features. | Use data-driven TensorFlow models for a wide range applications from music annotation to synthesis. | Find if a song is happy, sad, aggressive or relaxed. |
| **Key detection** | **Onset detection** | **Segmentation** | **Beat tracking** |
| Find a key of a music piece. | Detect onsets (and transients) in an audio signal. | Split audio into homogeneous segments that sound alike. | Estimate beat positions and tempo (BPM) of a song. |
| **Melody extraction** | **Audio fingerprinting** | **Cover song detection** | **Spectral analysis** |
| Estimate pitch in monophonic and polyphonic audio. | Extract fingerprints from any audio source using the Chromaprint algorithm. | Identify covers and different versions of the same music piece. | Analyze spectral shape of an audio signal. |
| **Loudness metering** | **Audio problems detection** | **Voice analysis** | **Synthesis** |
| Use various loudness meters | Identify possible audio quality | Voice activity detection and | Analyze, transform and |

# III. A Plug for nbdev

Python library development via
Jupyter notebooks

By Jeremy Howard & fast.ai crowd, but
doesn't require fastai

is "literate programming" instantiated:
code, docs & tests *are one*

Built-in CI via GitHub Actions

Hawley & Morrison (JASA-EL 2022)
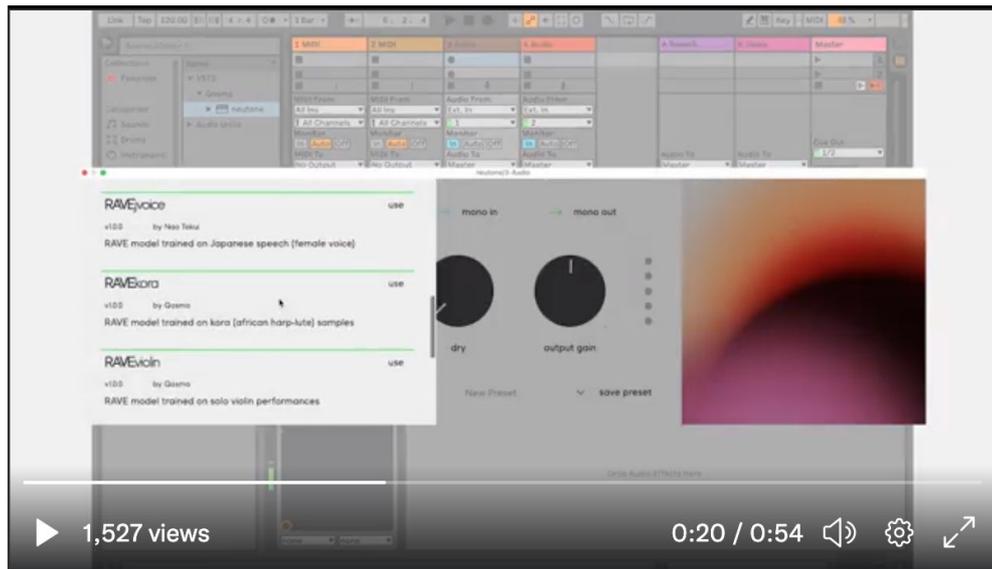found it ***essential for*** working
efficiently (**staying sane**)

# III. WE ALL SHOULD BUILD MORE



My attempt at interactive viz of PyTorch layer activations: images + oscilloscope.
To some student: *Please* fork this and make it *your own* & make it *good*.

# III. Post-Talk: New One! <u>Neutone</u> by Quosmo

"AI audio plugin & community, Bridging the gap between AI research and creativity"



Christian Steinmetz
@csteinmetz1

This is great for researchers and audio engineers! A platform for integrating your own deep learning models into plugins along with a simple way to share them. Excited to see where this goes.

**For audio creators**

Neutone makes AI technologies accessible for all to experiment with. You'll find transformative AI audio instruments that will spark endless creative possibilities.

**For AI researchers**

Neutone is a go-to platform for you to share real-time AI audio processing models with potential users in the audio production community.

1,527 views                                                0:20 / 0:54

Download Plugin                    Submit Your Model