



THE DEVELOPER'S
CONFERENCE

O que a IA pode extrair dos arquivos de **áudios** (**Audio insights**)

Diego Augusto Silva • Data Scientist • diegosilva@ciandt.com



[linkedin.com/in/diegoaugusto/](https://www.linkedin.com/in/diegoaugusto/)

Julho • 2019



The next **20 minutes**:

- Challenges
- Databases
- Applications

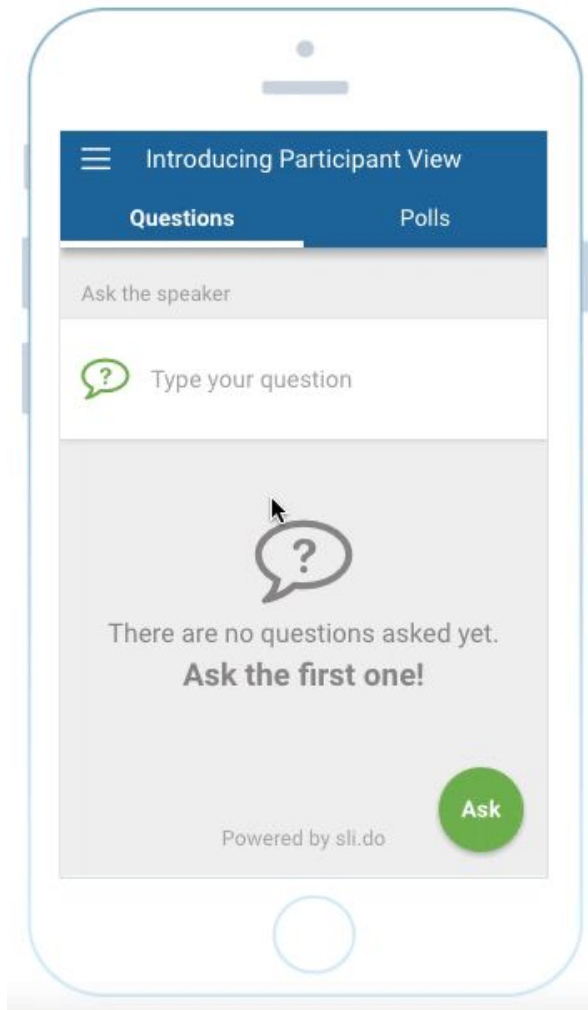


The next **20 minutes**:

- Challenges
- Databases
- Applications



Deep neural nets' result in
state of the art performance!



sli.do

Código: **iaudio**



Friendly reminder

Audio: Refers to an electrical signal (recorded).

Sound: When electrical signal is converted into audible acoustical pressure.

Speech/Voice: Sound produced by a person's.

Speech/Voice AI:

- Recognition (ASR)
- Synthesizer (TTS)
- Biometric
- Translation

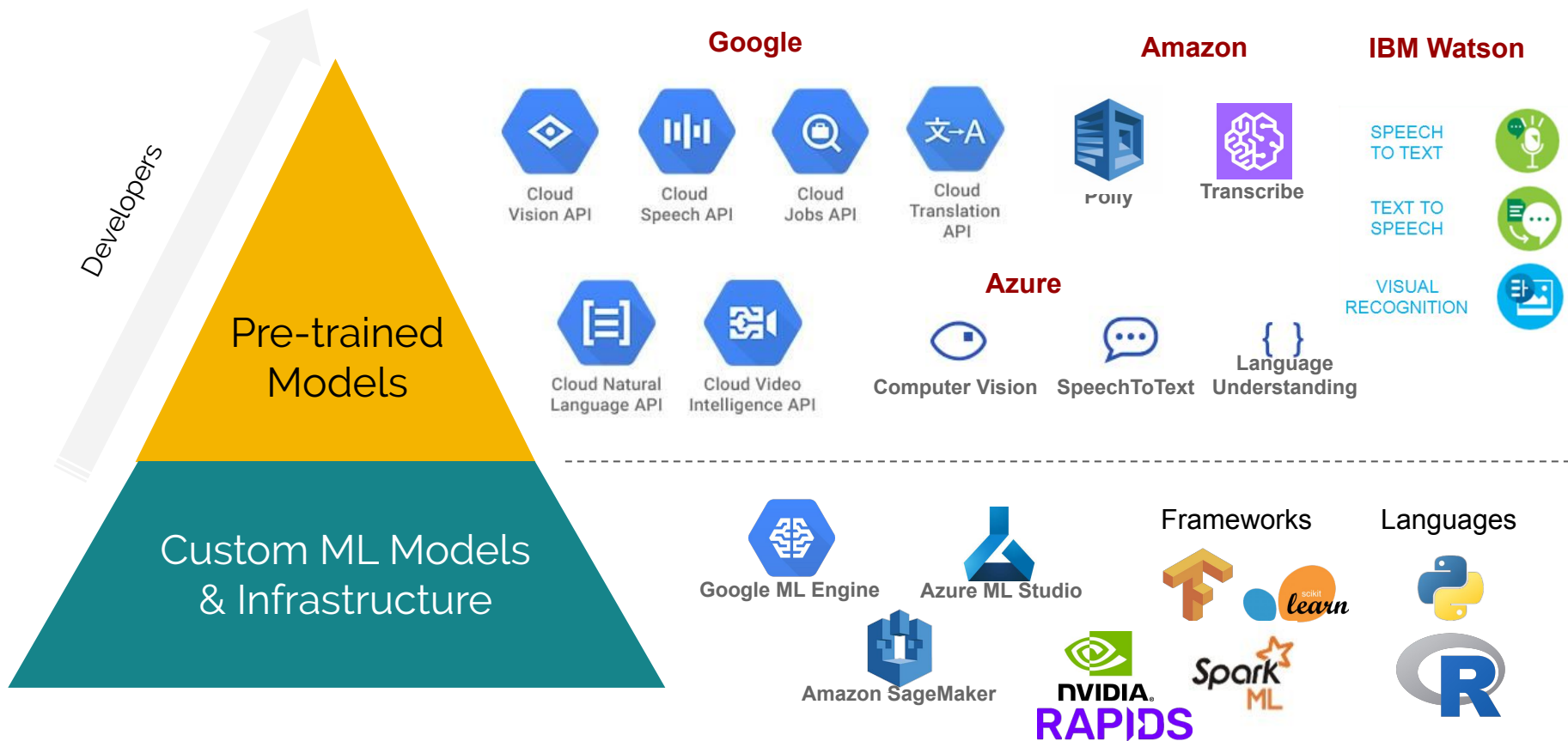
Music AI:

- Fingerprinting
- Genre
- Make

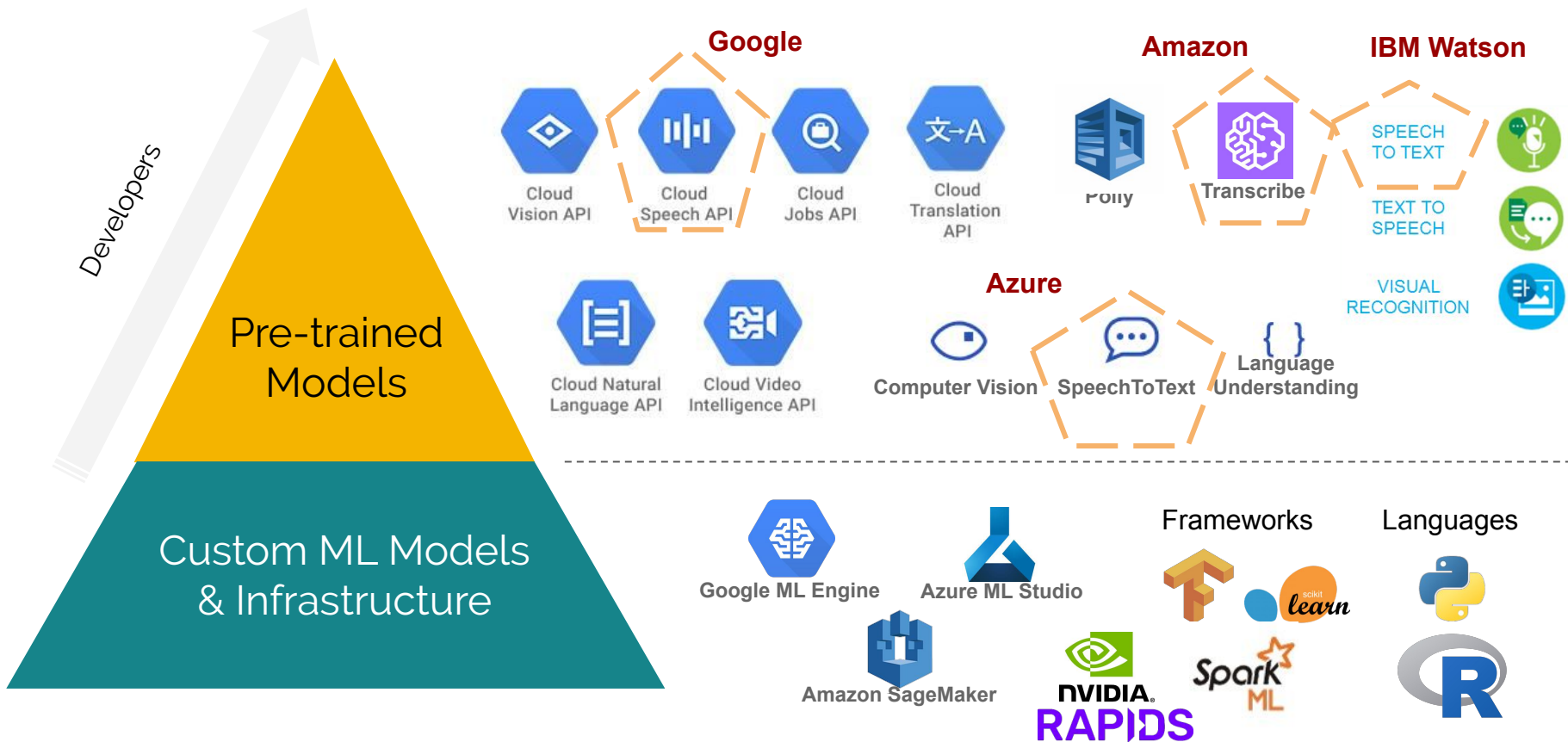
Sound AI:

- Scene Recog.
- Event Recog.

Framework and API



Framework and API



Why aren't sound events/scenes recognition done yet?

A successful **events/scenes recognition** system **requires**:

- The right representation.
- The right classification architecture.
- The right training data.



Why aren't sound events/scenes recognition done yet?

A successful **events/scenes recognition** system **requires**:

- The right representation.
- The right classification architecture.
- The right training data.

Discriminative model remains somewhat valid:

$$p(\textit{car alarm} \mid \textit{beep})$$

Why aren't sound events/scenes recognition done yet?

A successful **events/scenes recognition** system **requires**:

- The right representation.
- The right classification architecture.
- The right training data.

Discriminative model remains somewhat valid:

$$p(\text{car alarm} \mid \text{beep})$$

Oh no!!! It's a hard...

$$p(\text{car alarm} \mid \text{beep and a lots of garbage})$$

Audio Processing Challenges:

(1) Capture mode: Close-talking (headset) and Far-field (distant)

- IHM - Individual Headset Microphone(s)
- SDM - Single Distant Microphone
- MDM - Multiple Distant Microphones

(2) Distortion:

- Reverberation, Signal-to-Noise ratio, ...

(3) Events co-occurring:

- Music (jingle)
- Speech (overlap)

Audio Processing Challenges:

(1) Capture mode: Close-talking (headset) and Far-field (distant)

- IHM - Individual Headset Microphone(s)
- SDM - Single Distant Microphone
- MDM - Multiple Distant Microphones

(2) Distortion:

- Reverberation, Signal-to-Noise ratio, ...

(3) Events co-occurring:

- Music (jingle)
- Speech (overlap)
- **Noise (background)**

Audio Processing Challenges:

(1) Capture mode: Close-talking (headset) and **Far-field (distant)**

- IHM - Individual Headset Microphone(s)
- SDM - Single Distant Microphone
- **MDM - Multiple Distant Microphones**

(2) Distortion:

- Reverberation, Signal-to-Noise ratio, ...

(3) Events co-occurring:

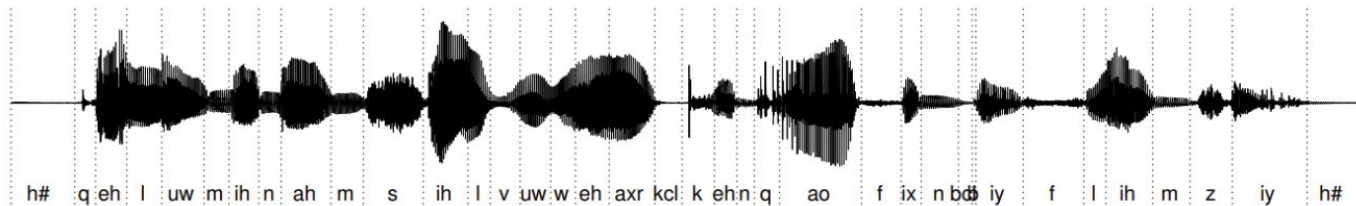
- Music (jingle)
- Speech (overlap)
- **Noise (background)**





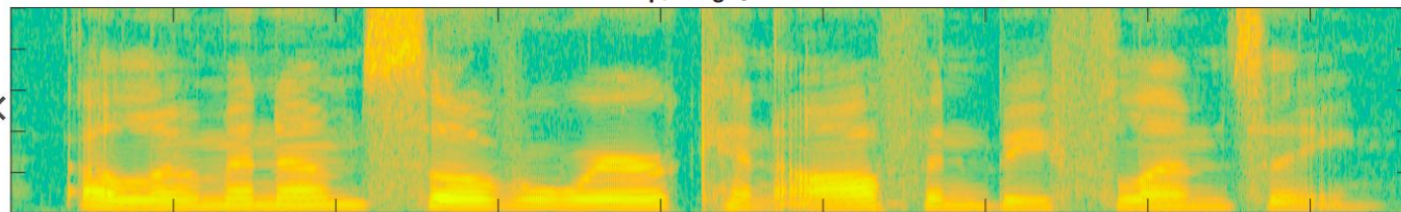
Audio representation

Raw Audio



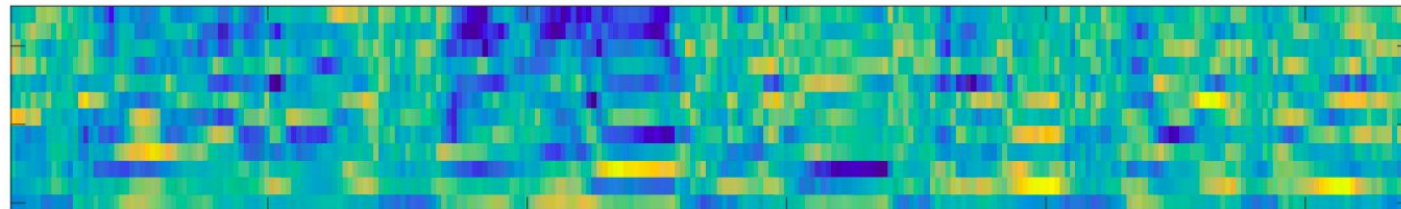
Raw Audio:
16k feats/sec.
shape: (1, 16000)

Spectrogram



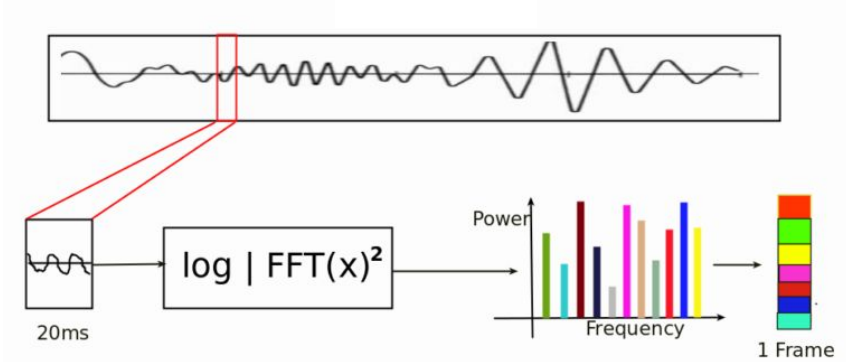
Spectrogram
~32k feats/sec.
shape: (128, 256)

MFCC Features



MFCC:
~3k feats/sec.
shape: (13, 256)

Audio representation



dog - 5-231762-A-0.wav



Databases: Labeled Acoustic Events

IEEE AASP Challenge on Detection/Classification of Acoustic Scenes/Events

- **DCASE** Challenges (since 2016): **Classification, Localization, Detec...**

Google Audioset:

- 527 sound events on YouTube videos (over 2 million recordings)

ESC-50:

- 50 classes (2k recordings)

Databases: Labeled Acoustic Events

IEEE AASP Challenge on Detection/Classification of Acoustic Scenes/Events

- **DCASE** Challenges (since 2016): **Classification, Localization, Detec...**

Google Audioset:

- 527 sound events on YouTube videos (over 2 million recordings)

ESC-50:

- 50 classes (2k recordings)

ESC-10:

- 10 classes - subset of ESC-50 (400 recordings)

Is it a Sound MNIST?
Yes, it is!



ESC-50 Public Results - Accuracy:

86.50%: CNN based + Features Fusion

- <https://pdfs.semanticscholar.org/f6fd/1be38a2d764d900b11b382a379efe88b3ed6.pdf>

84.90%: EnvNet-v2 + Data Augmentation

- <https://arxiv.org/pdf/1711.10282.pdf>

83.50%: CNN pretrained on AudioSet

- <https://arxiv.org/pdf/1711.01369.pdf>

IA vs Human?



ESC-50 Public Results - Accuracy:

86.50%: CNN based + Features Fusion

- <https://pdfs.semanticscholar.org/f6fd/1be38a2d764d900b11b382a379efe88b3ed6.pdf>

84.90%: EnvNet-v2 + Data Augmentation

- <https://arxiv.org/pdf/1711.10282.pdf>

83.50%: CNN pretrained on AudioSet

- <https://arxiv.org/pdf/1711.01369.pdf>

81.30%: Human listeners: Crowdsourcing experiment classification

- <http://karol.piczak.com/papers/Piczak2015-ESC-Dataset.pdf>

Applications:

Content-based Audio Search:

- Semantic search

Monitoring/Event Detection

- Secure: Gun shot, Window glass break, Door knock, Doorbell, Car alarm.
- Human: Shout, Laugh, Yawn, Cough, Snore.
- Nature: Rain, Wind.

Intelligent Content Processing (Recommendation)

- Dog bark, Cat meow, Bee buzz, Birds.
- Baby voice, Children playing.

Context Detection/Localization

- Home Indoor and Outdoor.
- Bus/Metro Station, Restaurant, Office.

Recap:

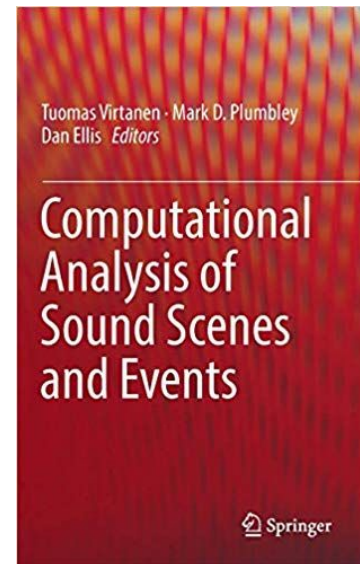
Sound is not speech!

- Not exactly the same problem as speech/music
- Diversity of sound production
- Context sensitive behavior

Computation cost matters.

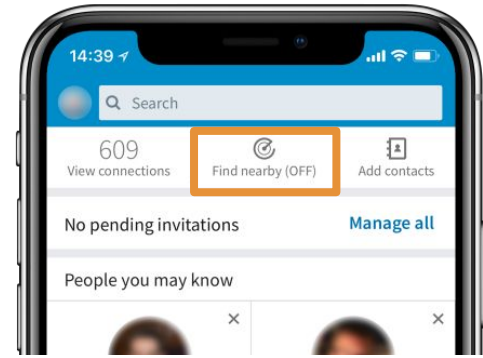
One against many.

Real world impact: Nonspeech Sound Event Captions





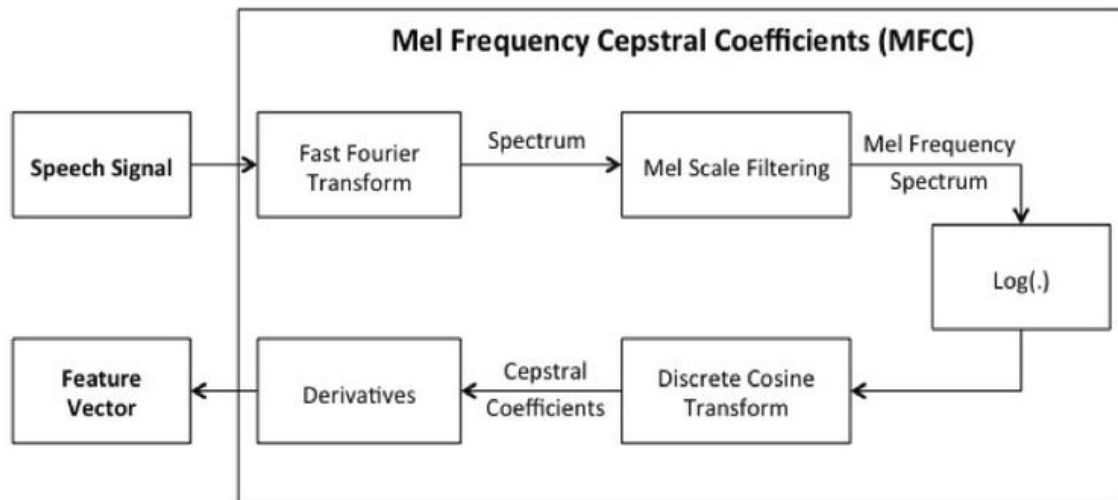
Diego Augusto Silva • Data Scientist • diegosilva@ciandt.com



linkedin.com/in/diegoaugusto

Audio Features

- MFCC [1, 2]
- ...



[1] <http://librosa.github.io/librosa/feature.html>

[2] http://kaldi-asr.org/doc/group_feat.html