



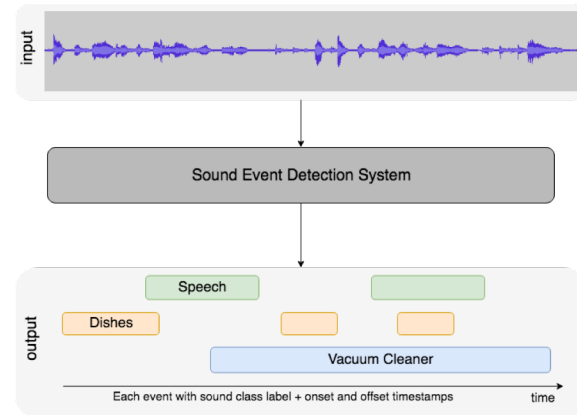
# Description and analysis of novelties introduced in DCASE Task 4 2022 on the baseline system

Francesca Ronchini<sup>1</sup>, Samuele Cornell<sup>2</sup>, Romain Serizel<sup>1</sup>, Nicolas Turpault<sup>1</sup>,  
Eduardo Fonseca<sup>3</sup>, Daniel P.W. Ellis<sup>3</sup>

<sup>1</sup>Université de Lorraine, CNRS, Inria, Loria <sup>2</sup>Università Politecnica delle Marche  
<sup>3</sup>Google, Inc.

# Why ?

- Necessity of real-world strongly annotated data ?
- Impact of external data and pre-trained models ?
- Environmental impact of our training models ?

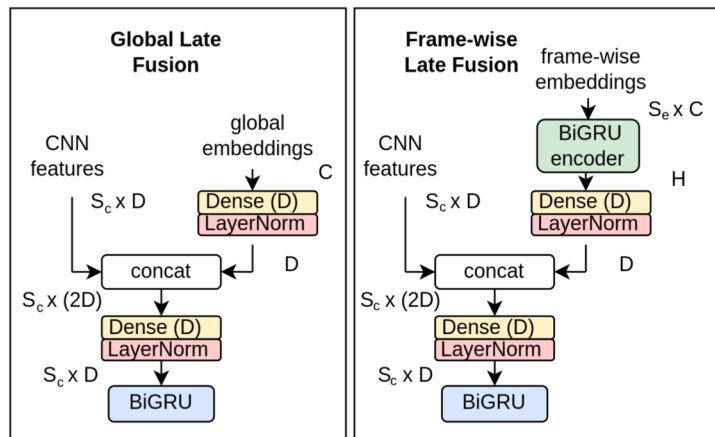


# Audioset strong !

- 3470 clips matching our labels
- These annotations are really expensive to get:
  - > Are they really helping ?
  - > Can't we get the same performance with synthetic data ?

# Embeddings from pre-trained models

- Using PANNs and AST:
  - > Global or frame embeddings ?
  
- Integrating embeddings is not trivial:
  - > Global fusion
  - > Frame-wise fusion
  - > Same problem with sound separation





# Energy-based metric ?

$$EW - PSDS = PSDS * \frac{kWh_{baseline}}{kWh_{submission}}$$

PSDS: polyphonic sound detection scores

$kWh_{baseline}$ : baseline energy consumption

$kWh_{submission}$ : system energy consumption

- How does this impact our models ?
- Are the systems using pre-trained models more efficient ?



*This is a simple suggestion, it should be improved → come & discuss !*

# Thank you

[nicolas.turpault@sonaide.fr](mailto:nicolas.turpault@sonaide.fr)