

Acoustic Space Learning: from Robots to Simulations

2009 - 2022

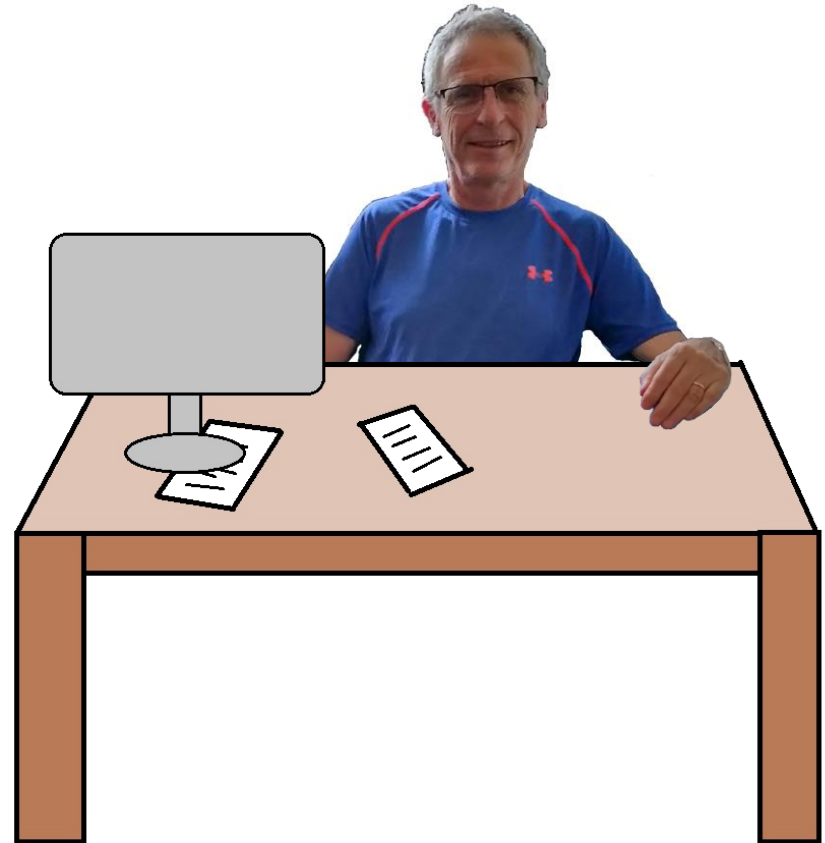


Antoine Deleforge, Inria Nancy – Grand Est, team MULTISPEECH

How I met Radu

2009

March 2009, Radu's office, Inria Grenoble

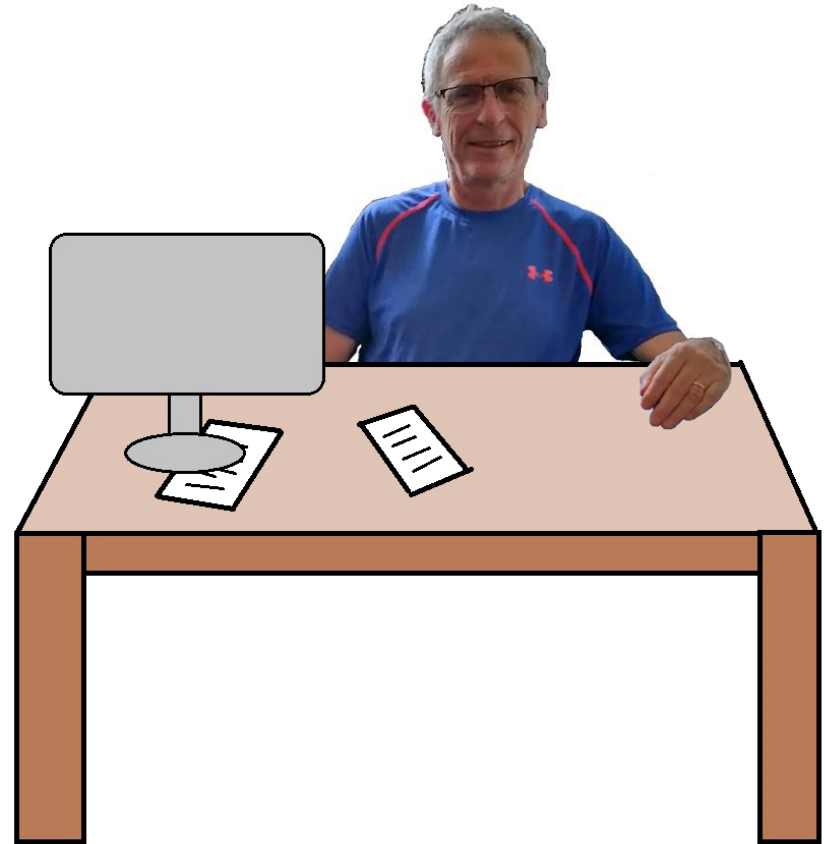


How I met Radu

2009

March 2009, Radu's office, Inria Grenoble

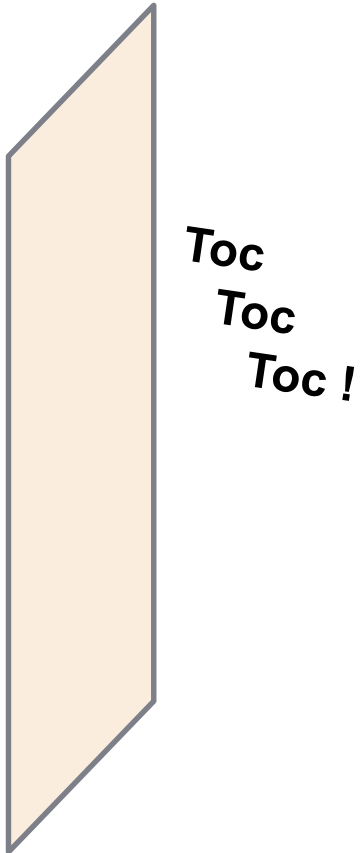
Toc
Toc
Toc!



How I met Radu

2009

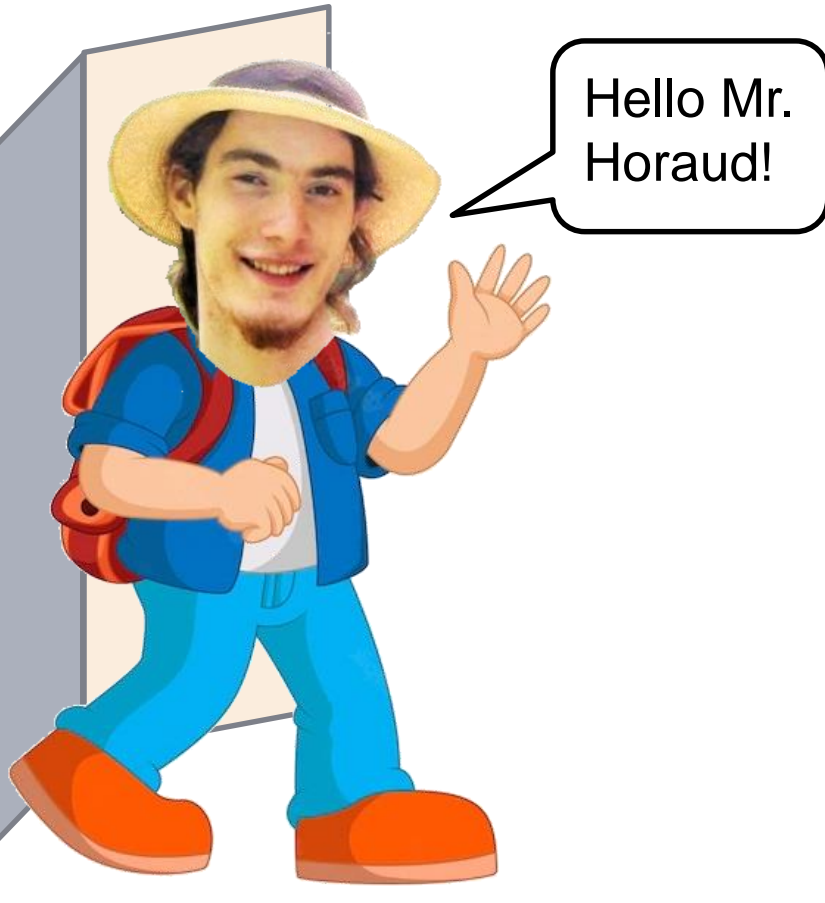
March 2009, Radu's office, Inria Grenoble



How I met Radu

2009

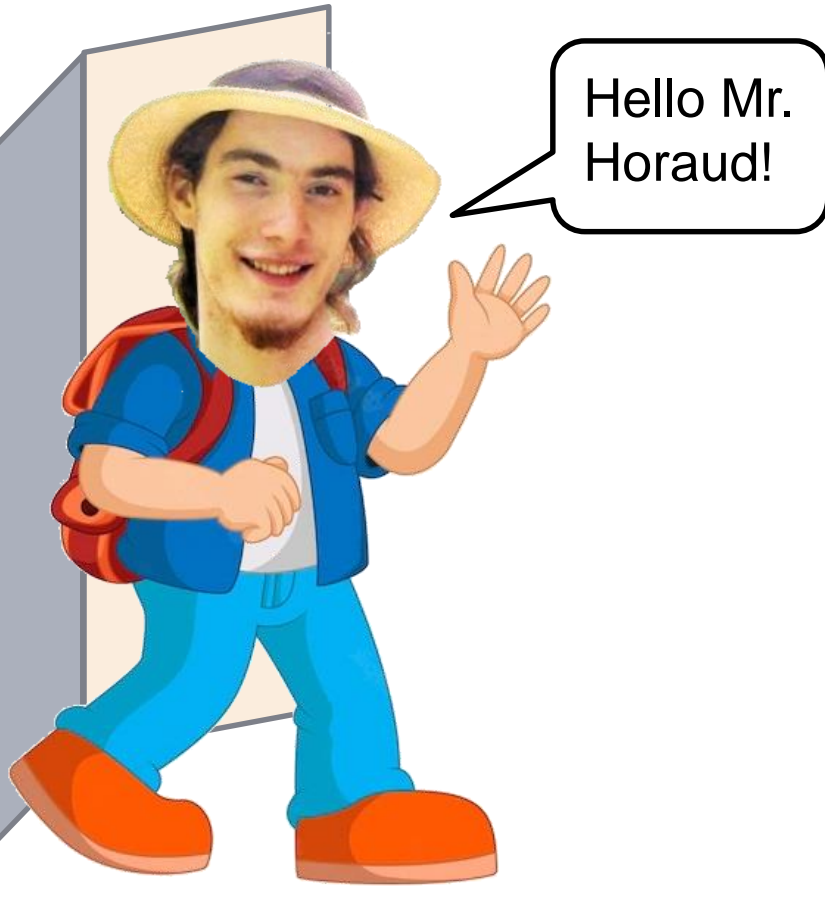
March 2009, Radu's office, Inria Grenoble



How I met Radu

2009

March 2009, Radu's office, Inria Grenoble



How I met Radu

March 2009, Radu's office, Inria Grenoble



I am a student at the engineering school ENSIMAG. I am pursuing a research master degree in **computer vision**, computer graphics & robotics.

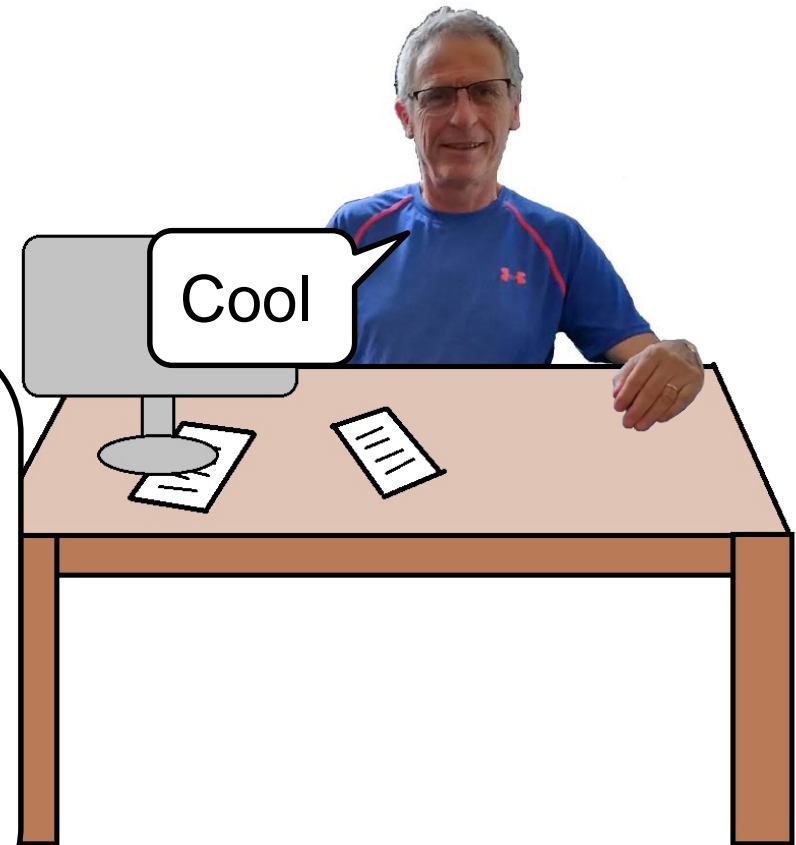


How I met Radu

March 2009, Radu's office, Inria Grenoble



I am a student at the engineering school ENSIMAG. I am pursuing a research master degree in **computer vision**, computer graphics & robotics.



How I met Radu

March 2009, Radu's office, Inria Grenoble



How I met Radu

March 2009, Radu's office, Inria Grenoble



You seem to be doing very interesting research in **computer vision** in your team!

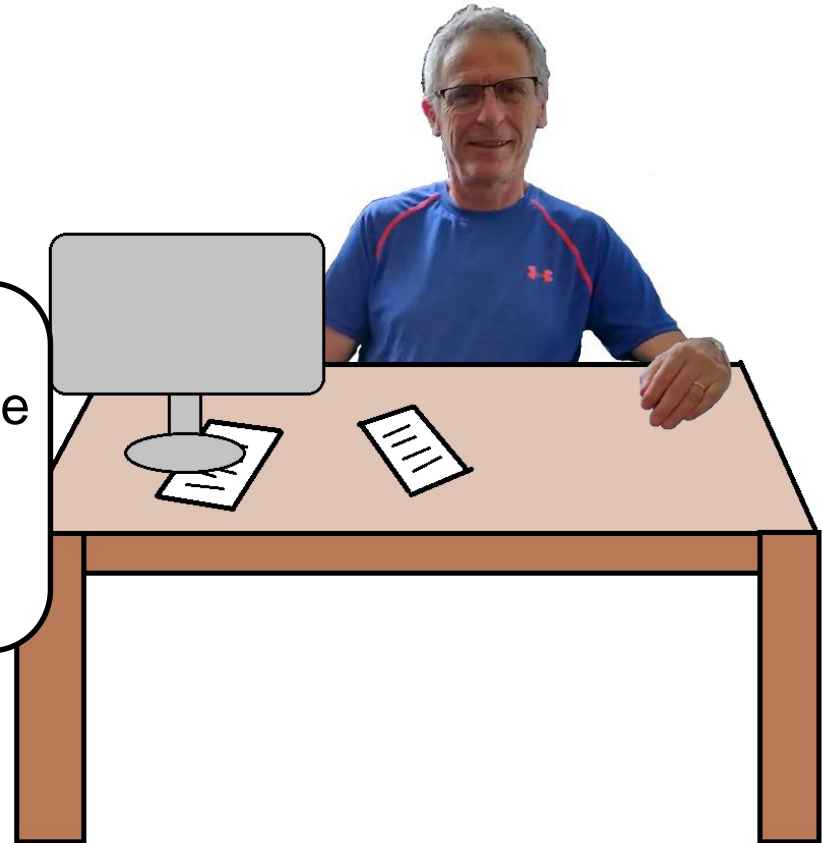
Indeed

How I met Radu

March 2009, Radu's office, Inria Grenoble

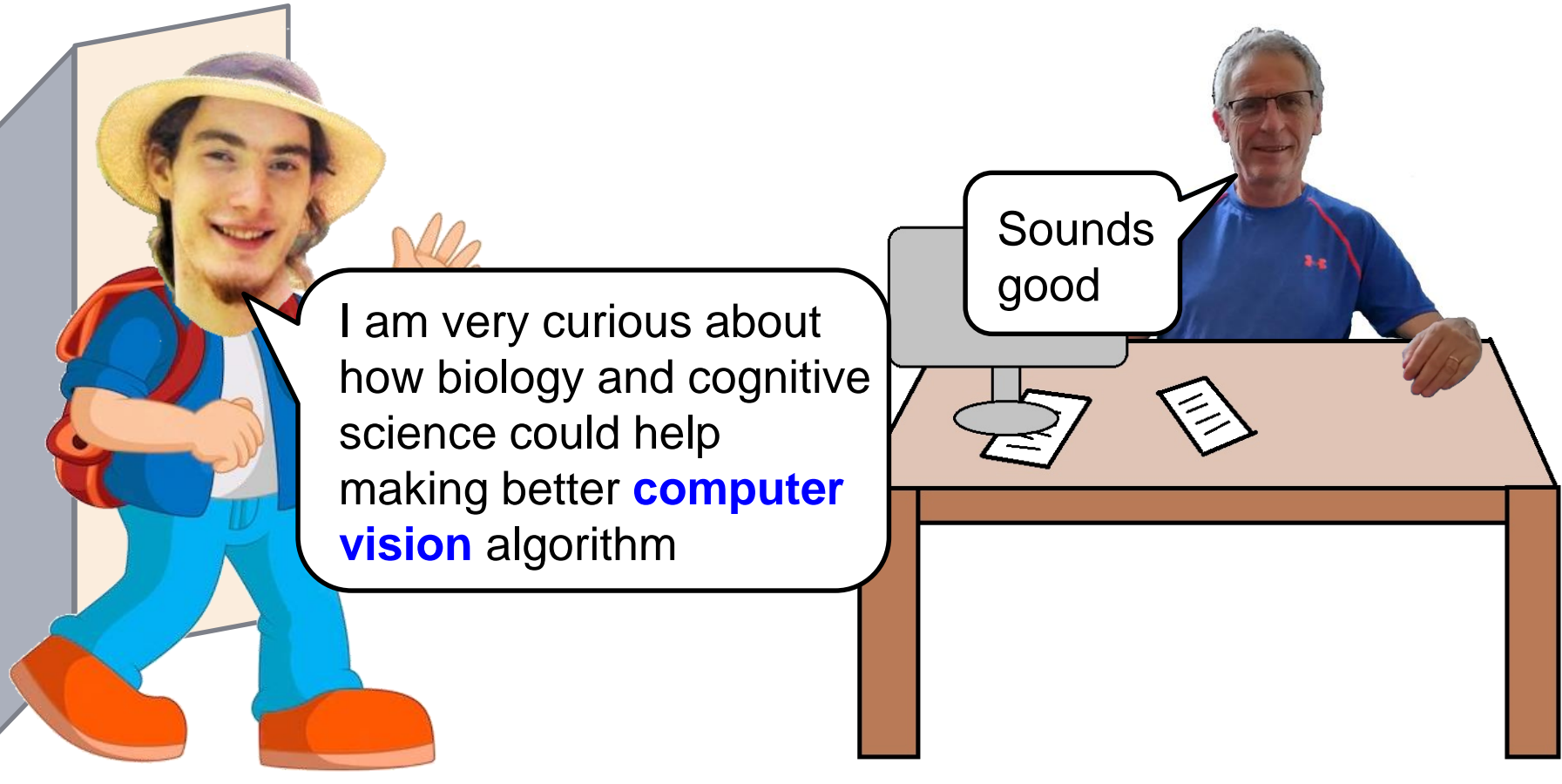


I am very curious about how biology and cognitive science could help making better **computer vision** algorithm



How I met Radu

March 2009, Radu's office, Inria Grenoble

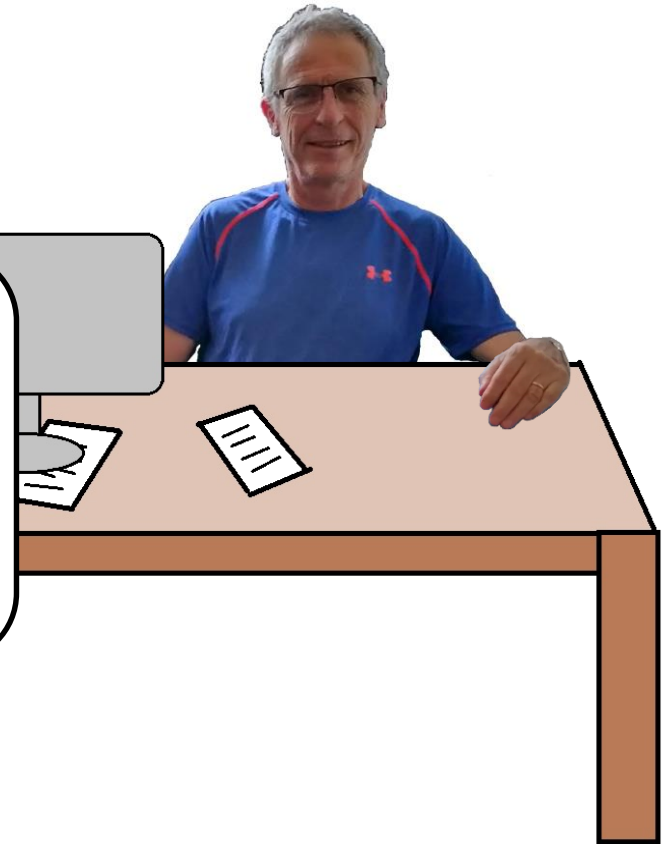


How I met Radu

March 2009, Radu's office, Inria Grenoble



It would be sooo nice to do an internship under your supervision, seeing how you know so much about **computer vision**, and all!



How I met Radu

March 2009, Radu's office, Inria Grenoble



It would be sooo nice to do an internship under your supervision, seeing how you know so much about **computer vision**, and all!

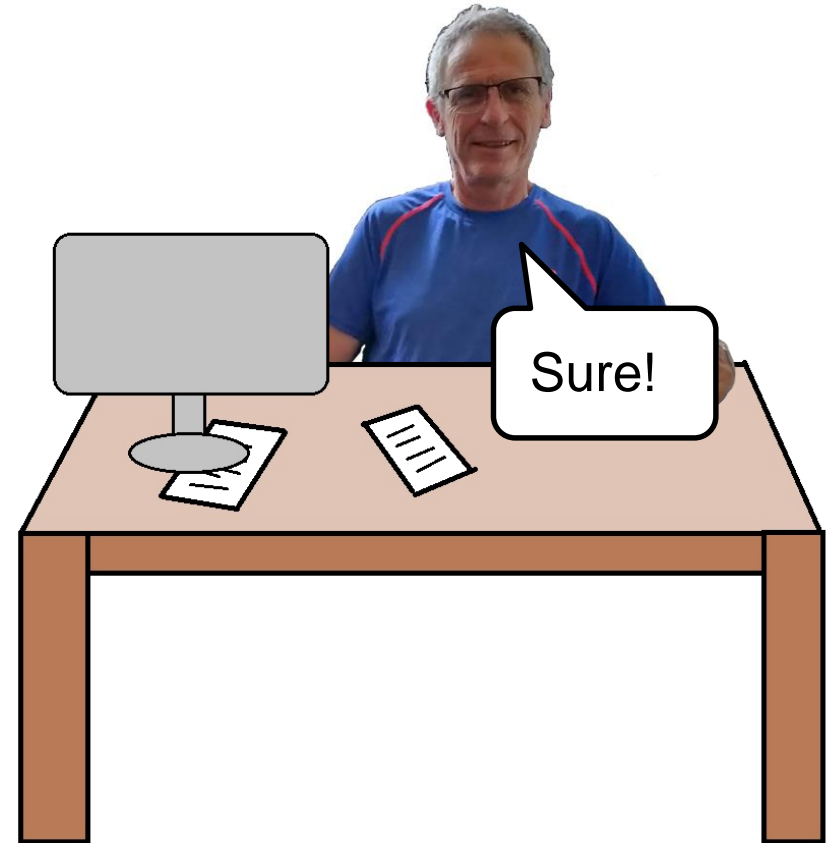


Sure!

How I met Radu

2009

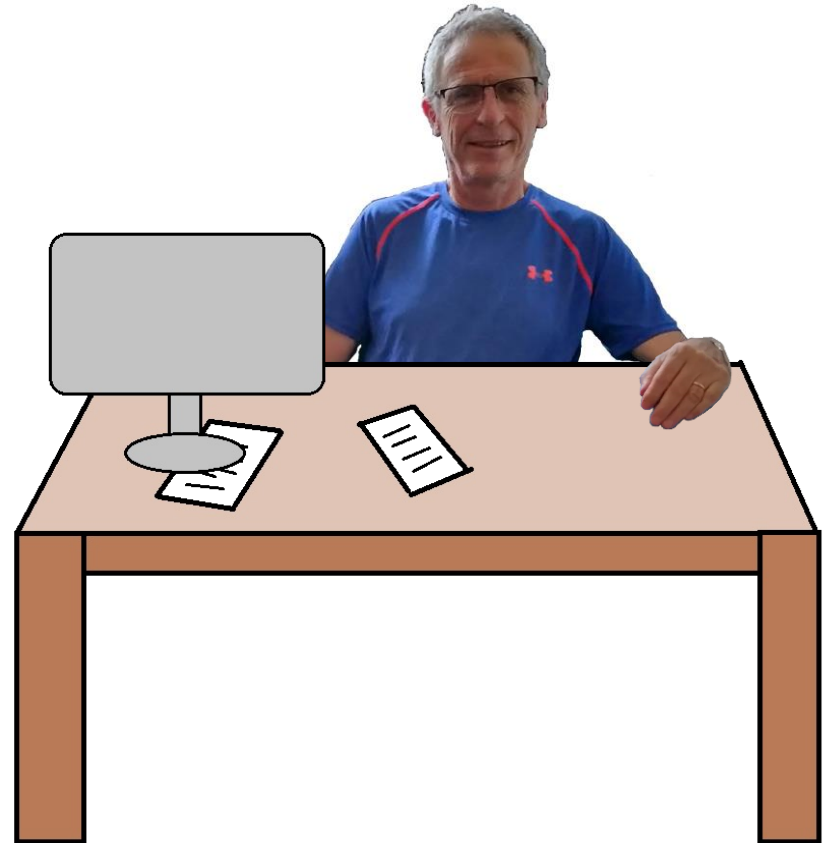
March 2009, Radu's office, Inria Grenoble



How I met Radu

2009

March 2009, Radu's office, Inria Grenoble



How I met Radu

2009

March 2009, Radu's office, Inria Grenoble



How I met Radu

2009

March 2009, Radu's office, Inria Grenoble



How I met Radu

2009

March 2009, Radu's office, Inria Grenoble



How I met Radu

2009

March 2009, Radu's office, Inria Grenoble



How I met Radu

2009

March 2009, Radu's office, Inria Grenoble



How I met Radu

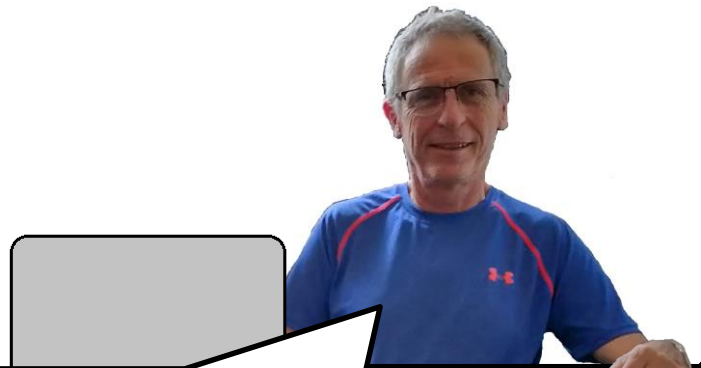
2009

March 2009, Radu's office, Inria Grenoble



How I met Radu

March 2009, Radu's office, Inria Grenoble

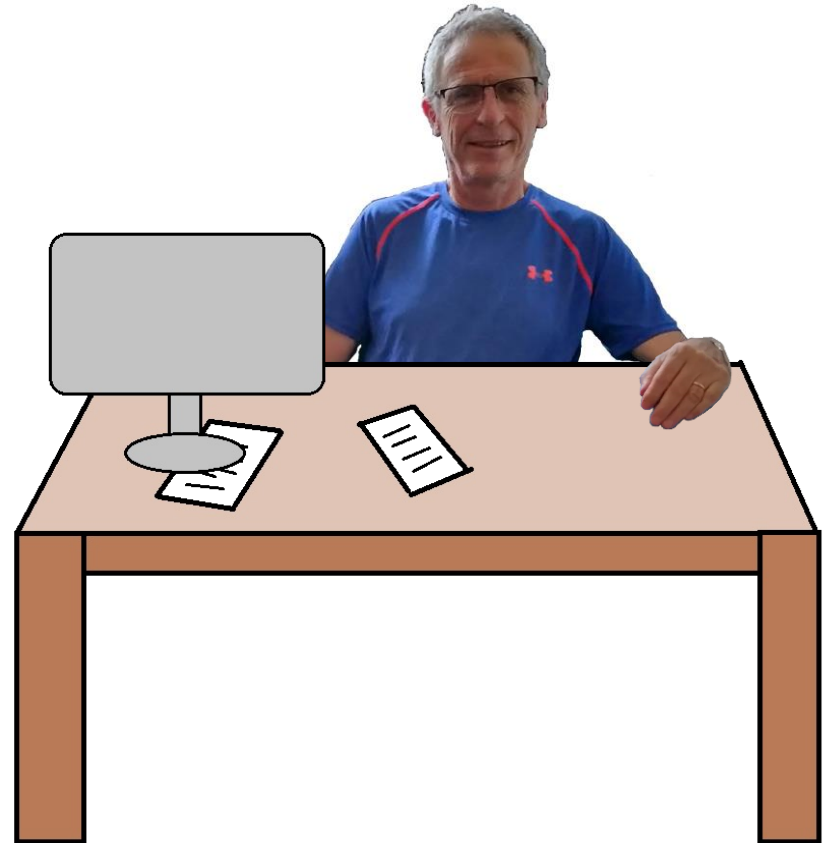


Something like:
Sound Source
Localization with a Robot:
Can We Use Biology?

How I met Radu

2009

March 2009, Radu's office, Inria Grenoble



How I met Radu

2009

March 2009, Radu's office, Inria Grenoble



How I met Radu

2009

March 2009, Radu's office, Inria Grenoble



How I met Radu

2009

March 2009, Radu's office, Inria Grenoble



How I met Radu

2009

March 2009, Radu's office, Inria Grenoble



See you next week!



ARTICLE  Communicated by J. Kevin O'Regan

A Sensorimotor Approach to Sound Localization

Murat Aytekin

aytekin@umd.edu

Neuroscience and Cognitive Science Program, University of Maryland, College Park, MD 20742, U.S.A.

Cynthia F. Moss

cmoss@psyc.umd.edu

Neuroscience and Cognitive Science Program, Department of Psychology and Institute of Systems Research, University of Maryland, College Park, MD 20742, U.S.A.

Jonathan Z. Simon

jzsimon@umd.edu

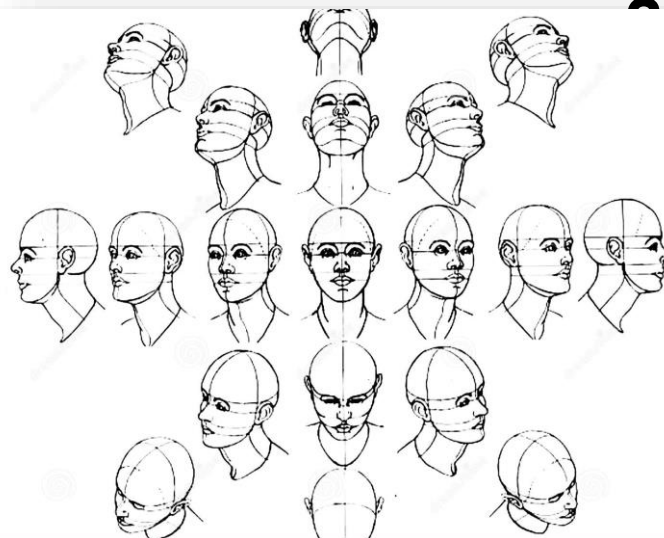
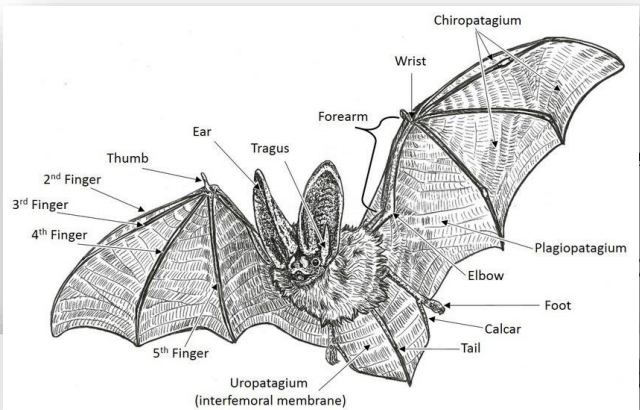
Neuroscience and Cognitive Science Program, Department of Electrical and Computer Engineering, Department of Biology, University of Maryland, College Park, MD 20742, U.S.A.

ARTICLE

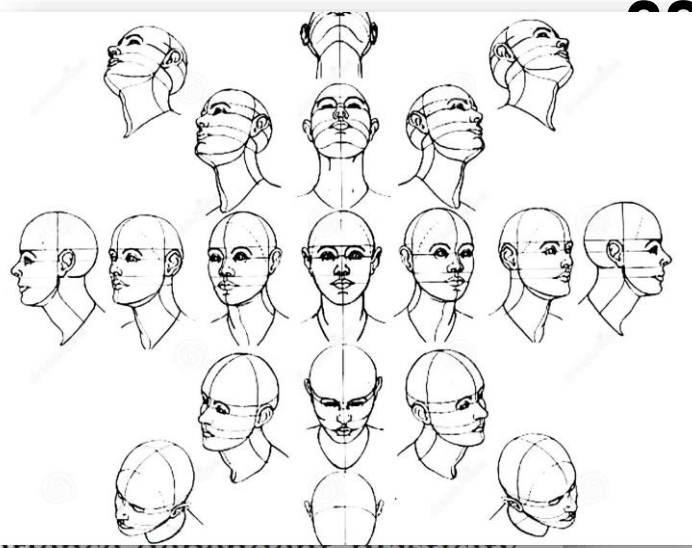
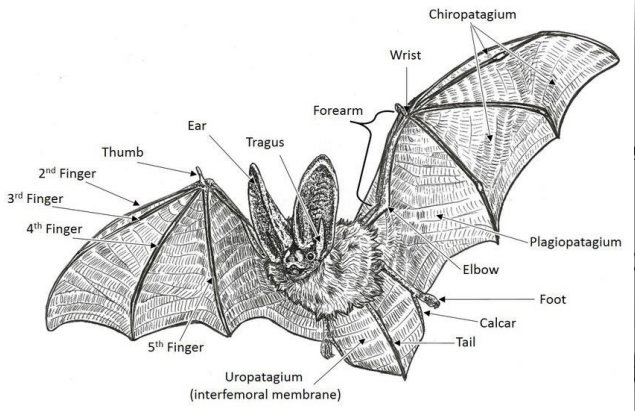
 Communicated by J. Kevin O'Regan

A Sensorimotor Approach to Sound Localization

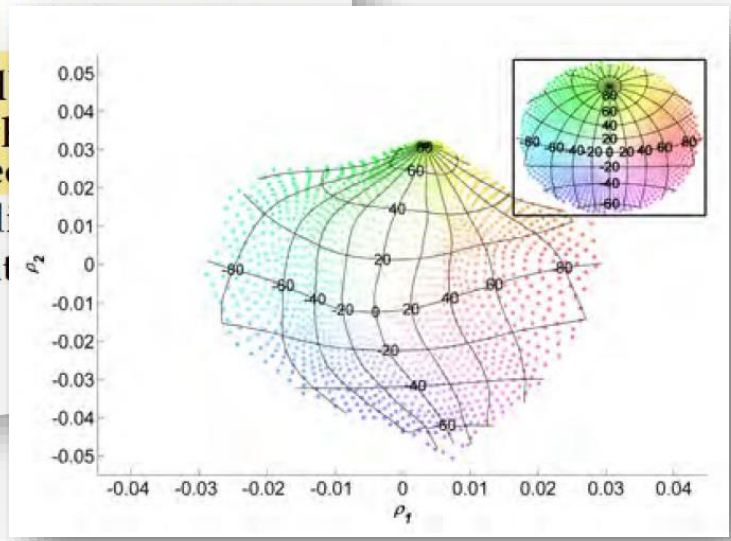
Sound localization is known to be a complex phenomenon, combining multisensory information processing, experience-dependent plasticity, and movement. Here we present a sensorimotor model that addresses the question of how an organism could learn to localize sound sources without any a priori neural representation of its head-related transfer function or prior experience with auditory spatial information. We demonstrate quantitatively that the experience of the sensory consequences of its voluntary motor actions allows an organism to learn the spatial location of any sound source. Using examples from humans and echolocating bats, our model shows that a naive organism can learn the auditory space based solely on acoustic inputs and their relation to motor states.



multisensory information processing, experience-dependent plasticity, and movement. Here we present a sensorimotor model that addresses the question of how an organism could learn to localize sound sources without any a priori neural representation of its head-related transfer function or prior experience with auditory spatial information. We demonstrate quantitatively that the experience of the sensory consequences of its voluntary motor actions allows an organism to learn the spatial location of any sound source. Using examples from humans and echolocating bats, our model shows that a naive organism can learn the auditory space based solely on acoustic inputs and their relation to motor states.

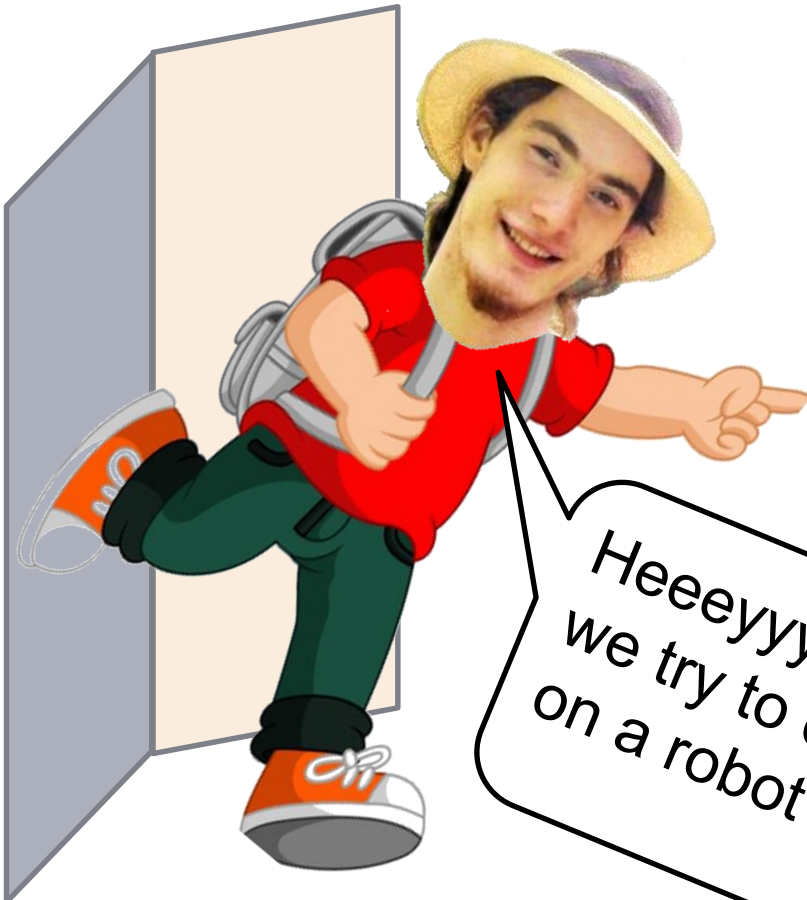


multisensory information processing, experience-dependent plasticity, and movement. Here we present a sensorimotor model that addresses the question of how an organism could learn to localize sound sources without any a priori neural representation of its head-related transfer function or prior experience with auditory spatial information. We show quantitatively that the experience of the sensory consequences of voluntary motor actions allows an organism to learn the spatial location of any sound source. Using examples from humans and elephants, our model shows that a naive organism can learn the auditory localization solely on acoustic inputs and their relation to motor states.



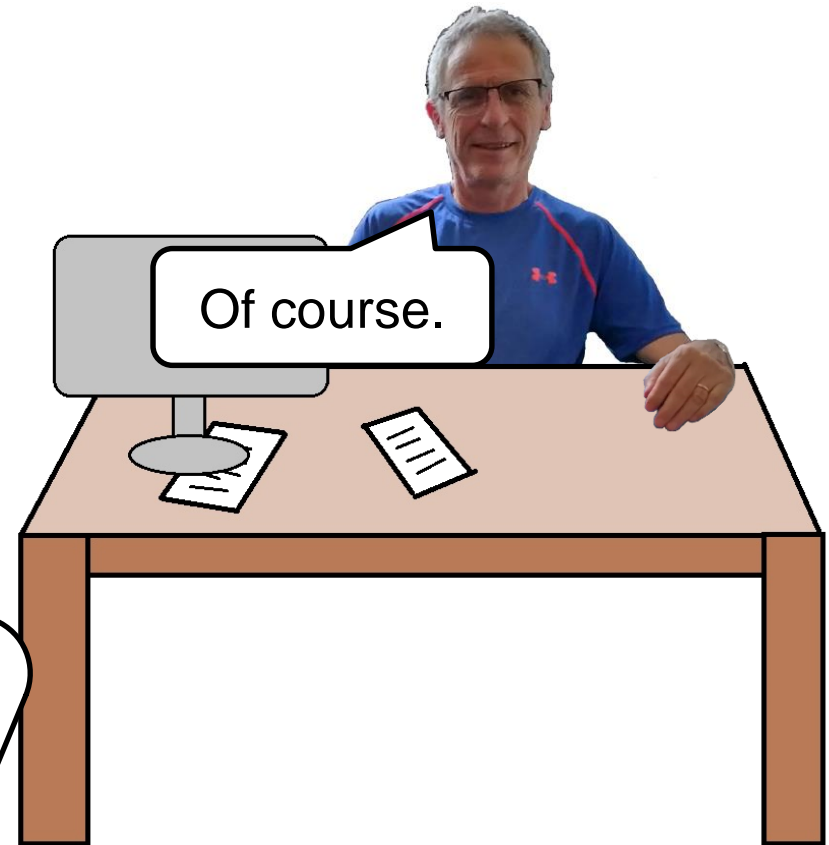
Toc
Toc
Toc!



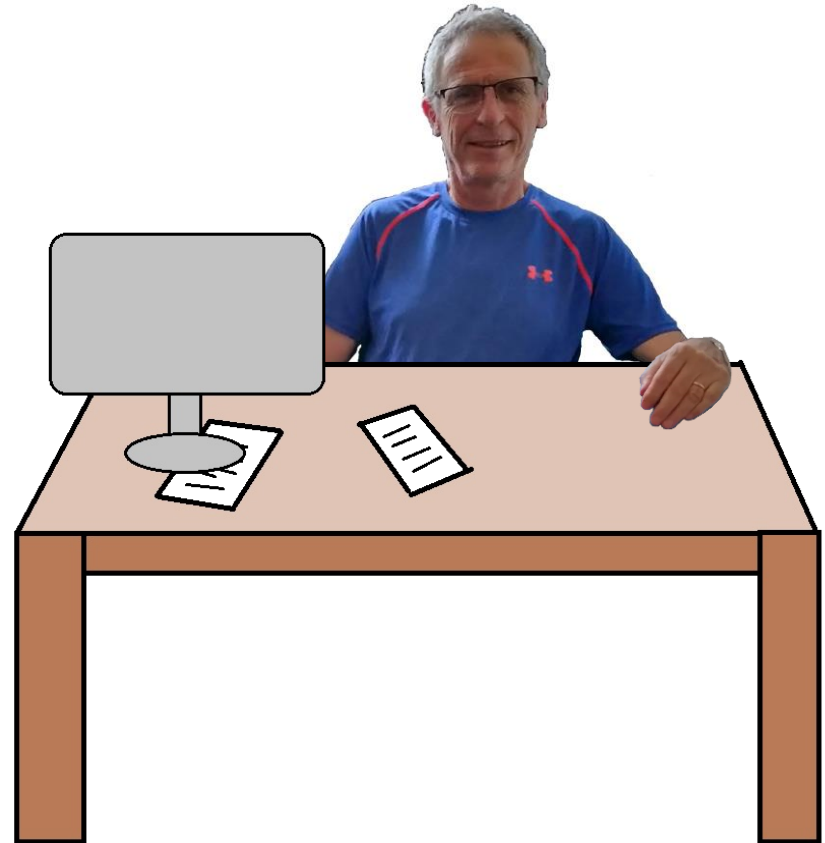


Heeeyyy! Can we try to do that on a robot??!!









The POPEYE Robot

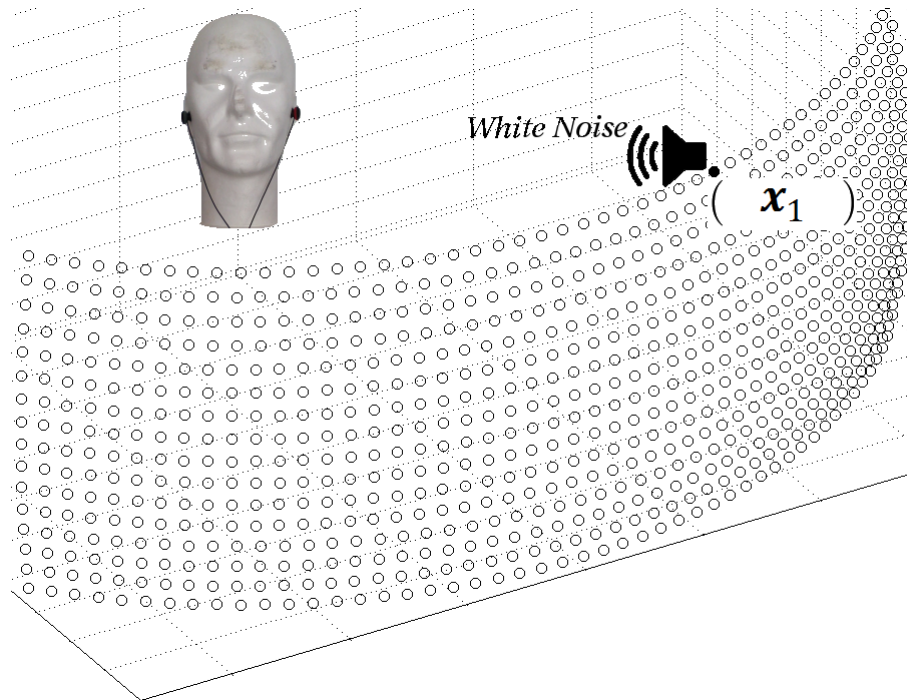
- Use a **white noise emitter** to obtain interaural cues in all frequencies
- Two methods: the **listener** is moved or the **emitter** is moved

Audio-motor sampling

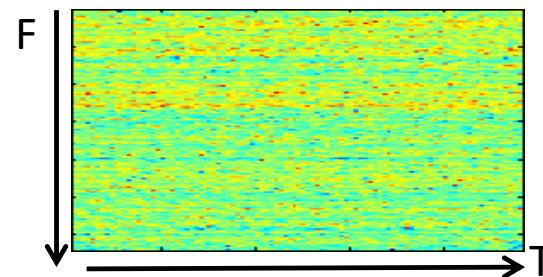
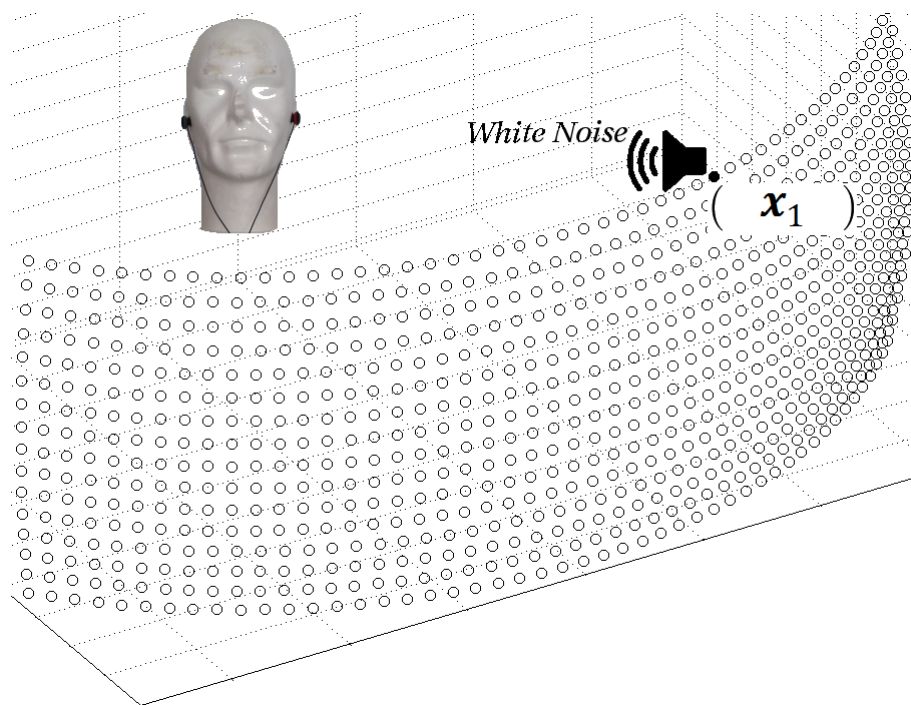
- 10,800 motor states (2 angles)
- 360° azimuth, 120° elevation

Audio-visual sampling

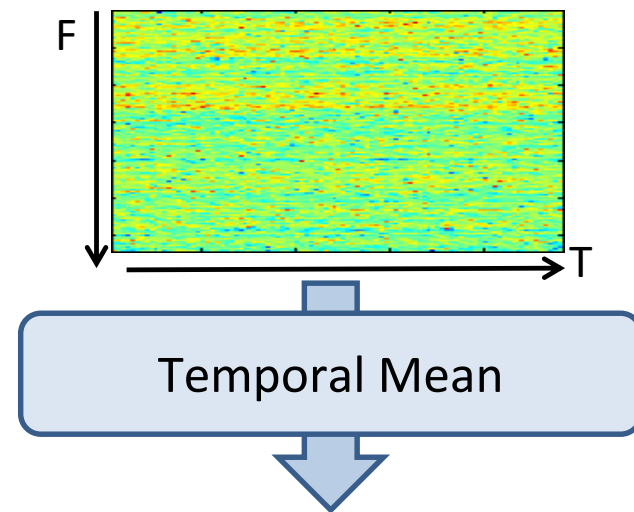
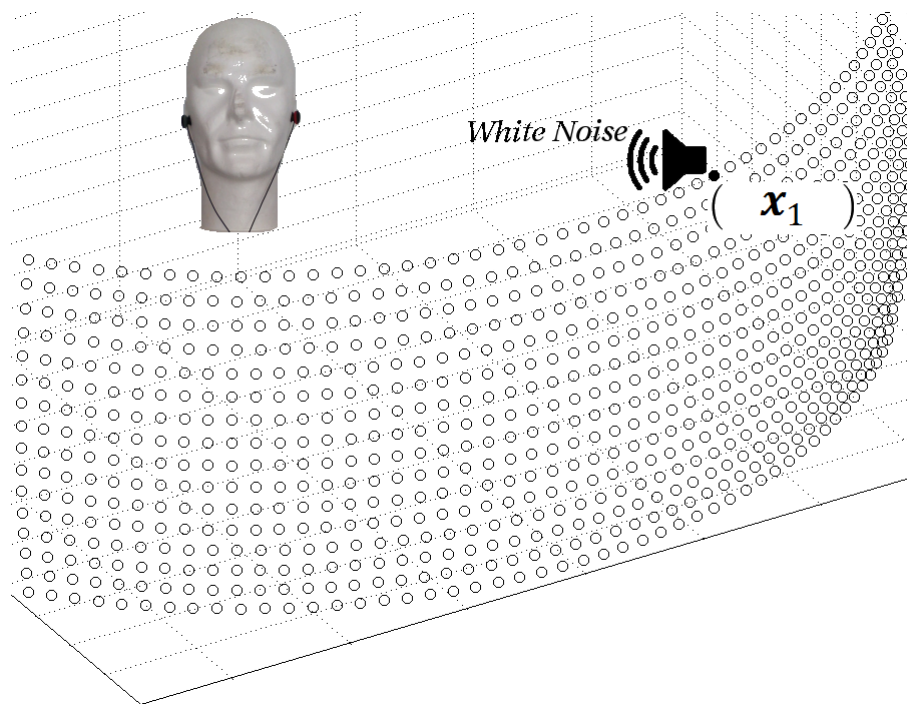
- 432 image positions (in pixels)
- Covering the camera field of view



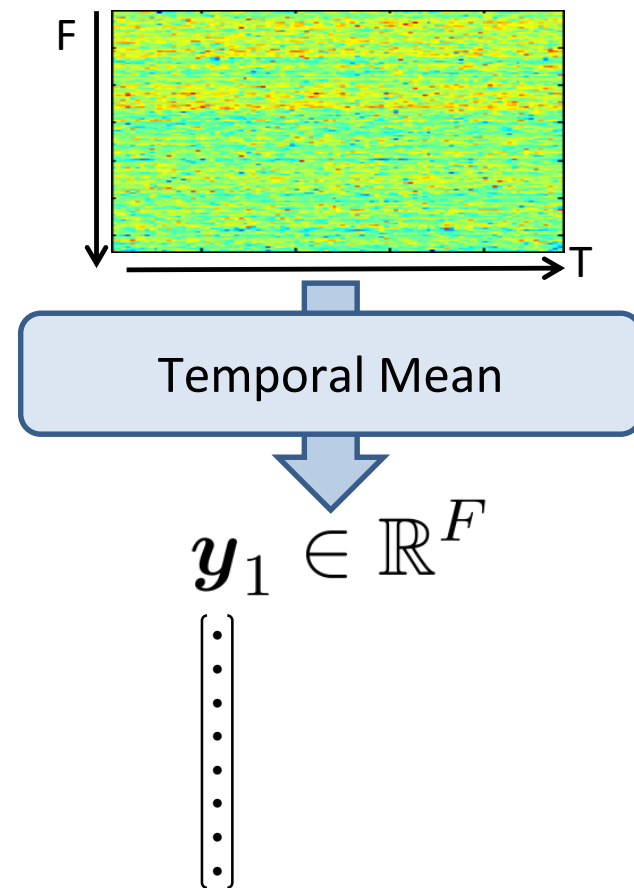
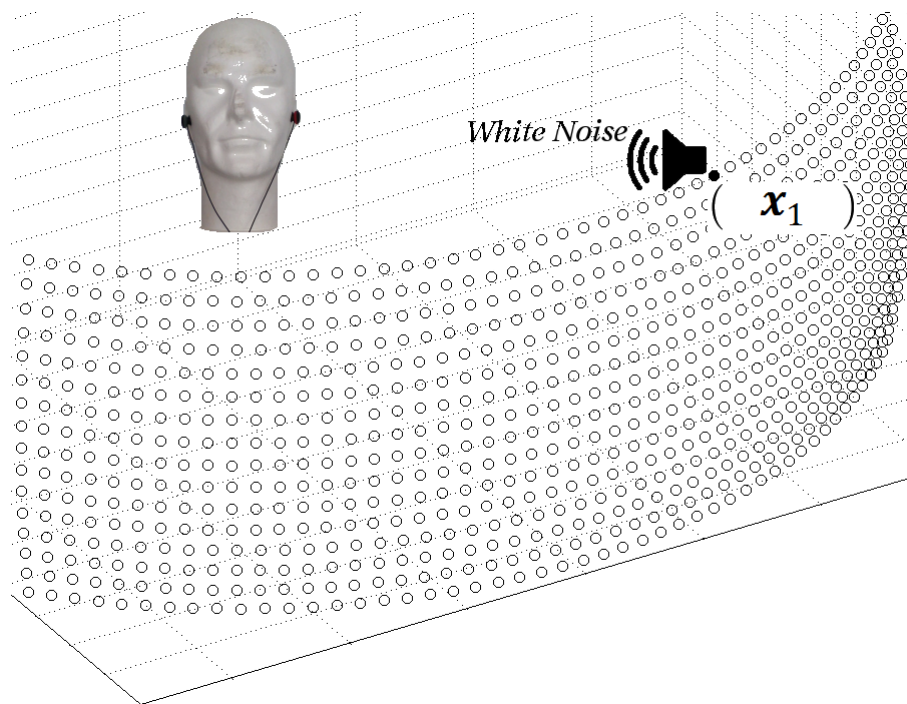
- Interaural Level Difference Spectrogram



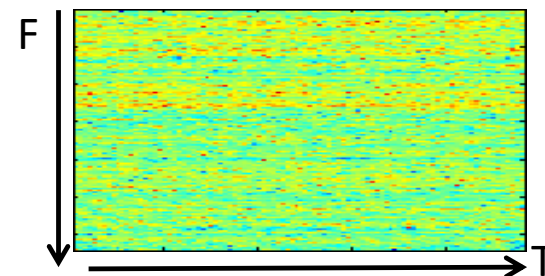
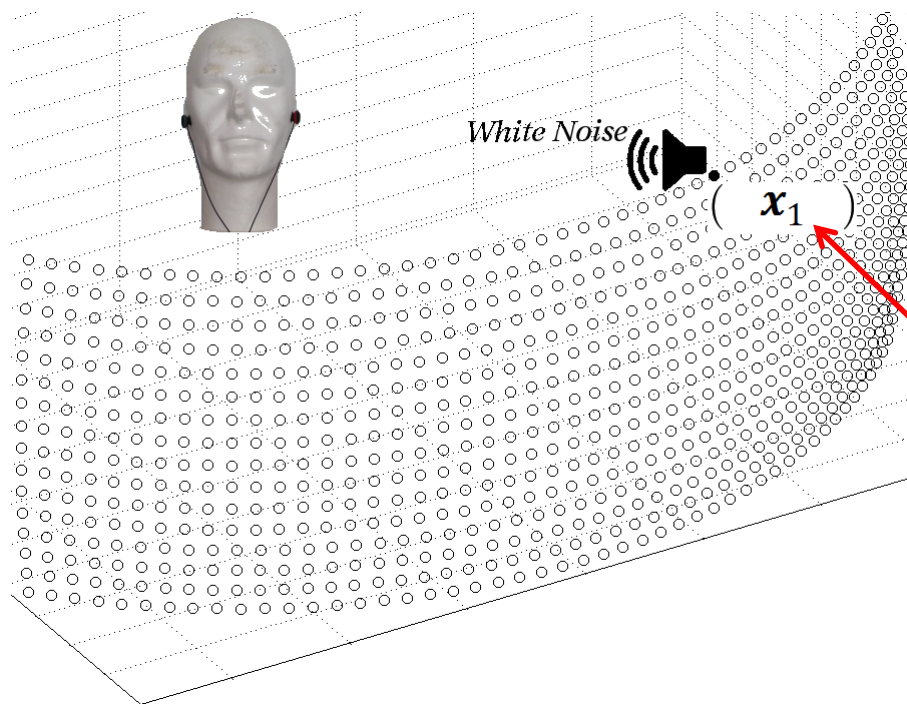
- Interaural Level Difference Spectrogram



- Interaural Level Difference Spectrogram

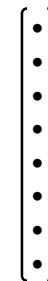


- Interaural Level Difference Spectrogram

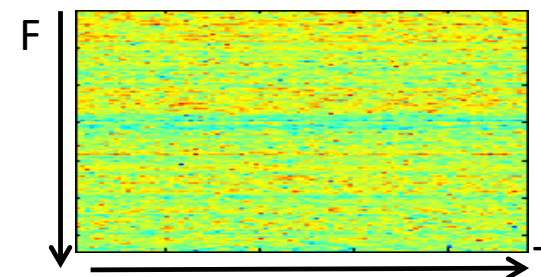
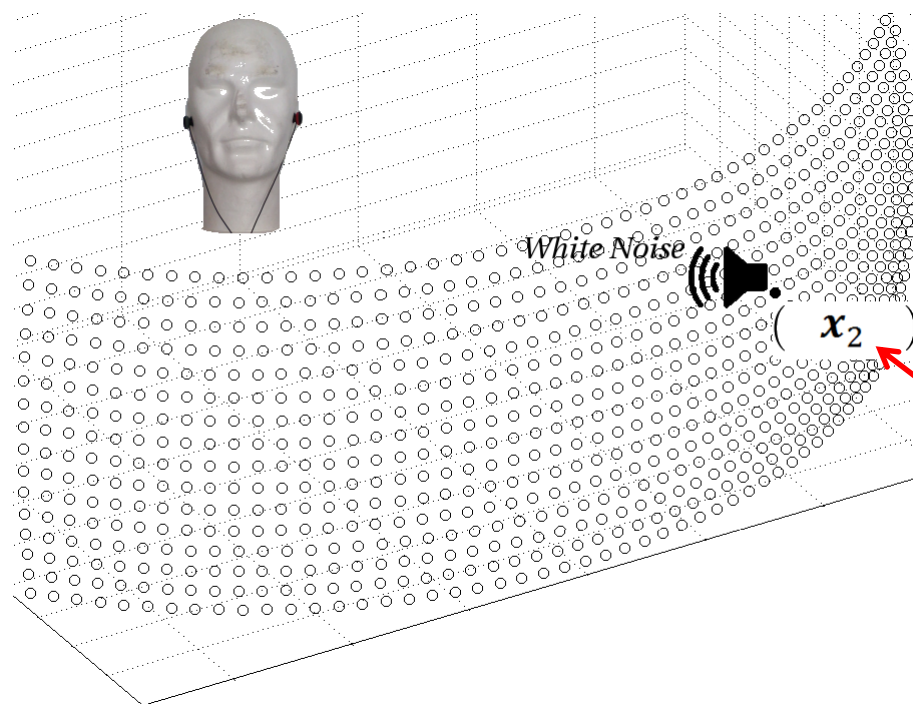


Temporal Mean

$$y_1 \in \mathbb{R}^F$$

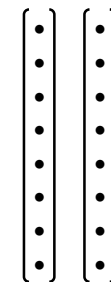


- Interaural Level Difference Spectrogram

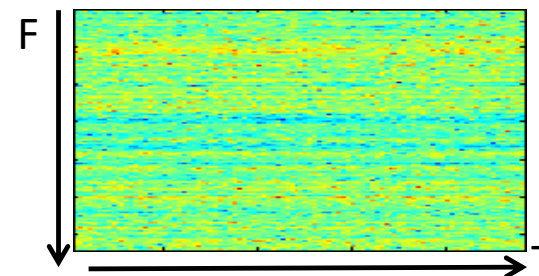
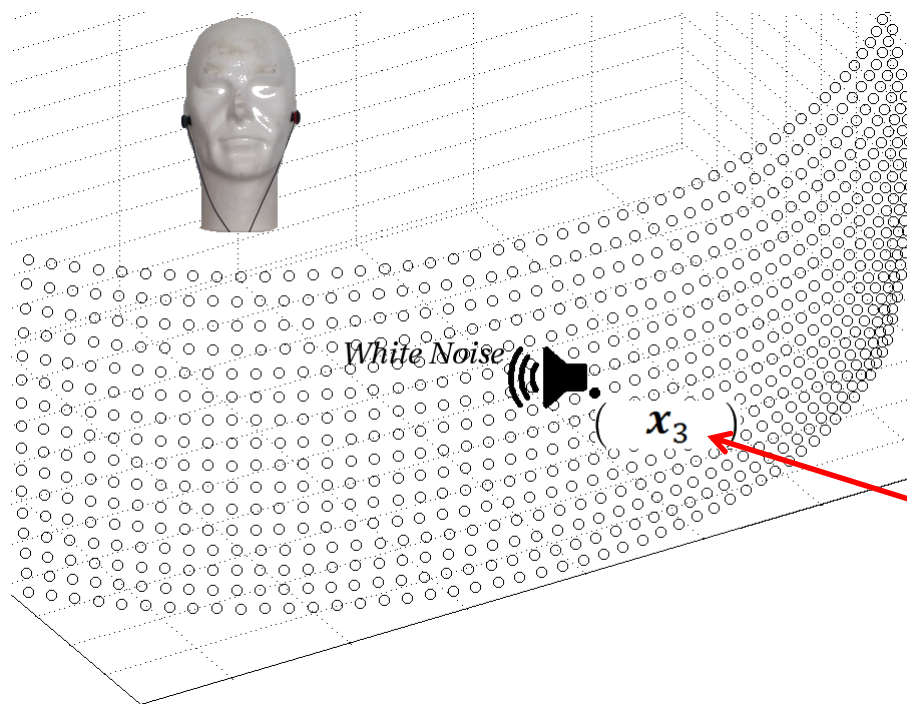


Temporal Mean

$$\mathbf{y}_2 \in \mathbb{R}^F$$

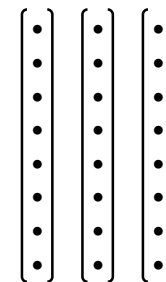


- Interaural Level Difference Spectrogram

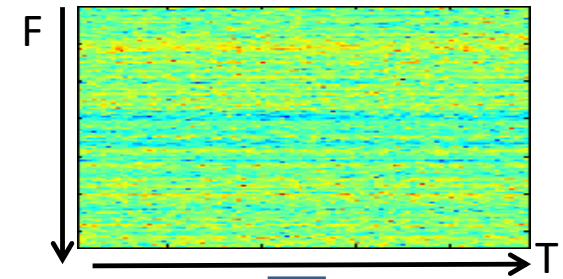
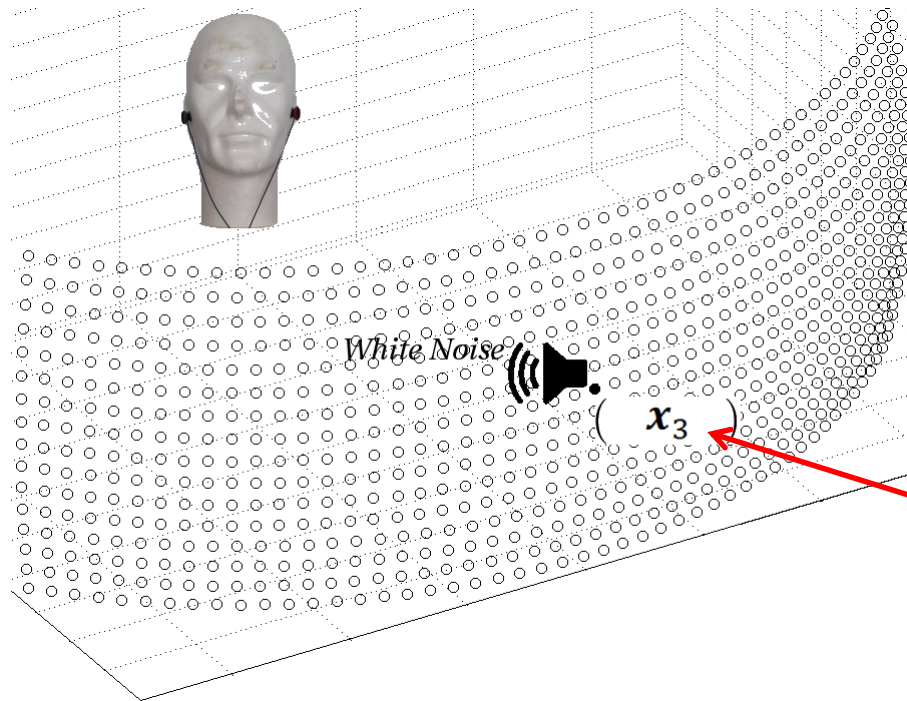


Temporal Mean

$$y_3 \in \mathbb{R}^F$$

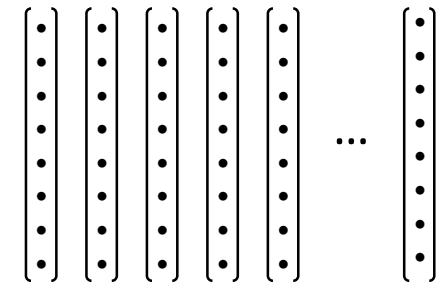


- Interaural Level Difference Spectrogram



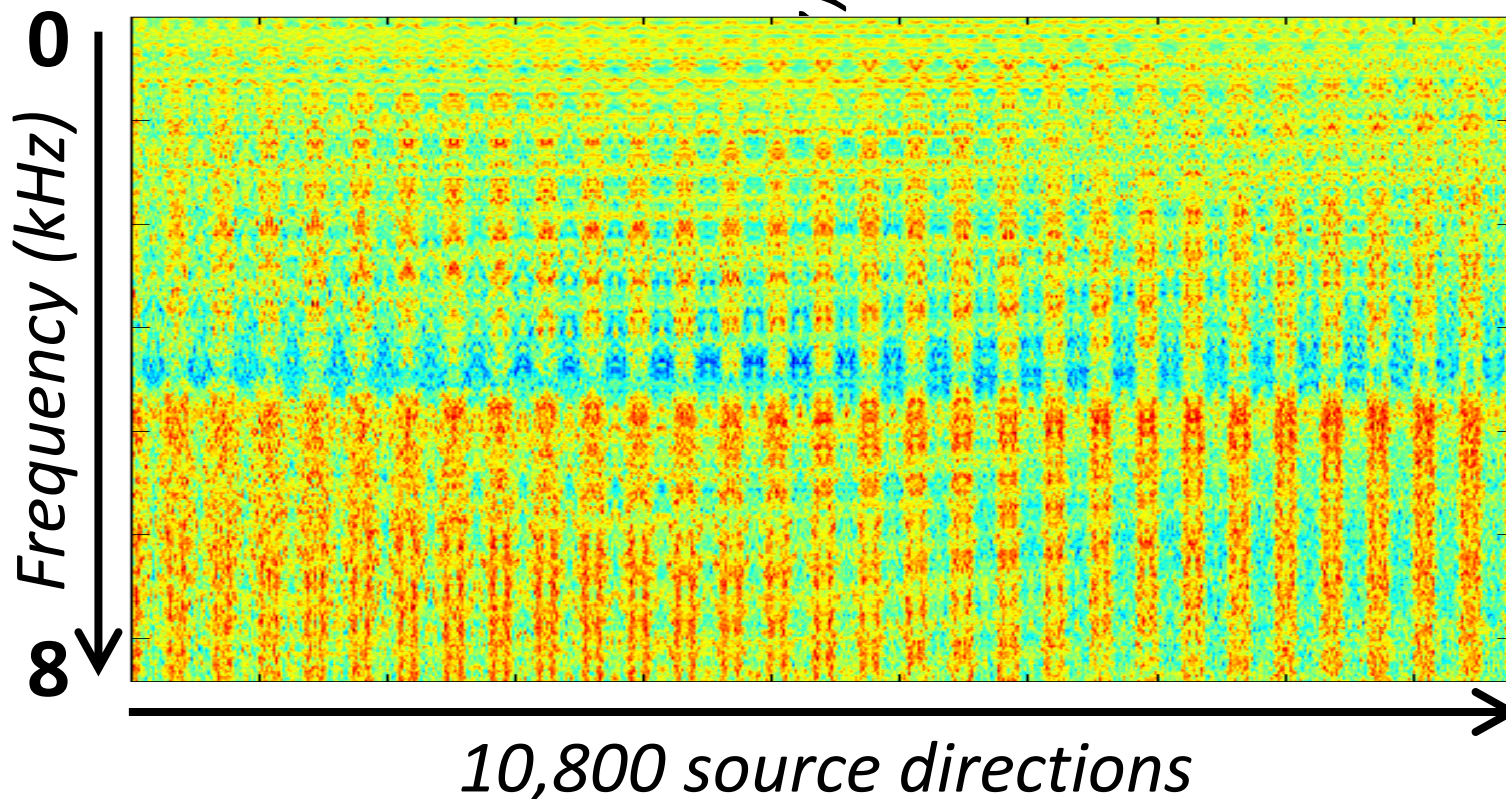
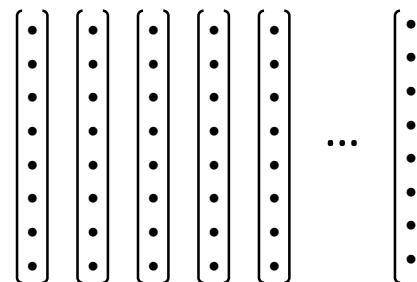
Temporal Mean

$$y_3 \in \mathbb{R}^F$$

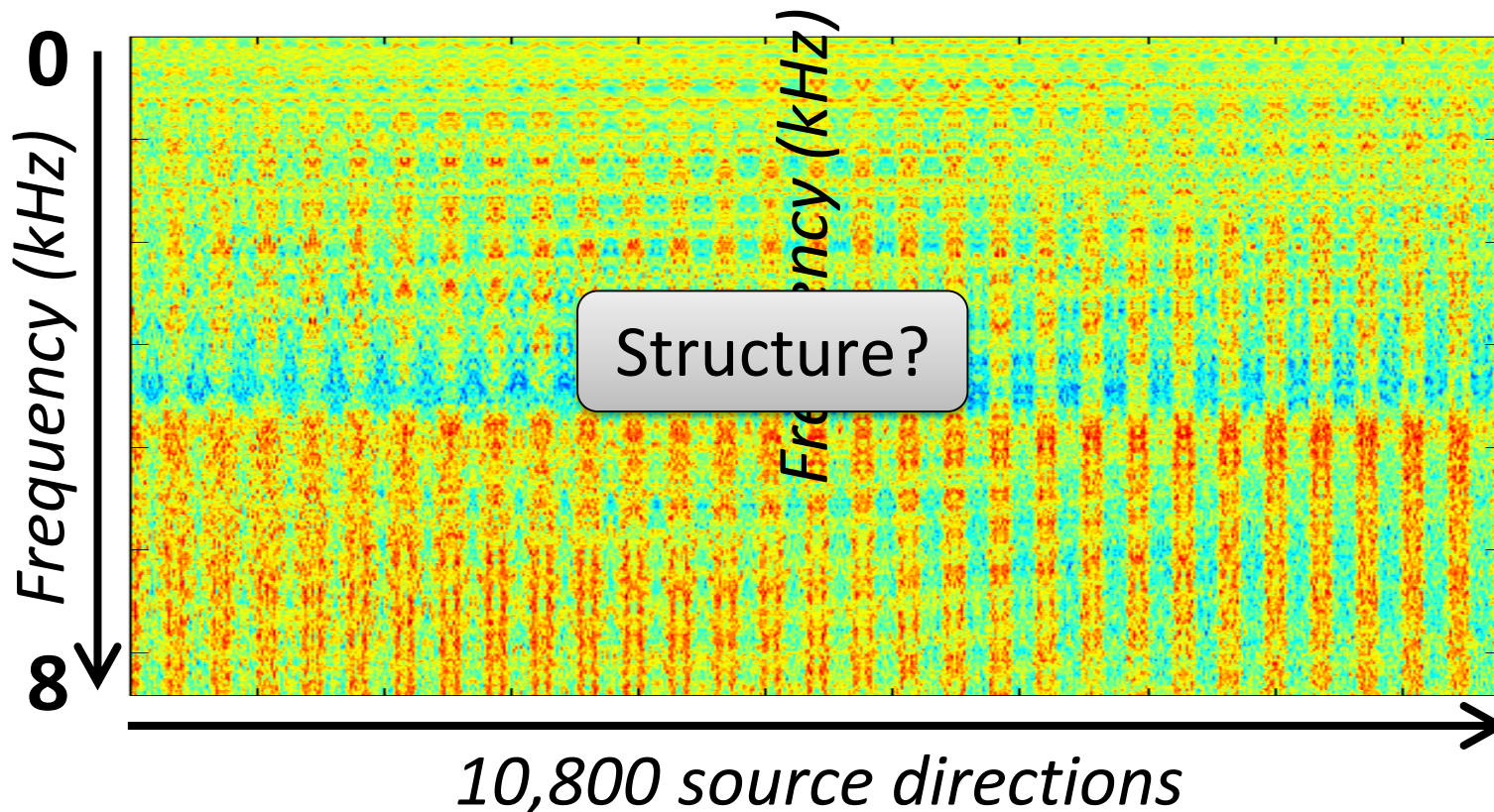
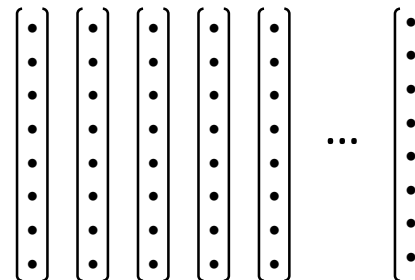


Acoustic Space

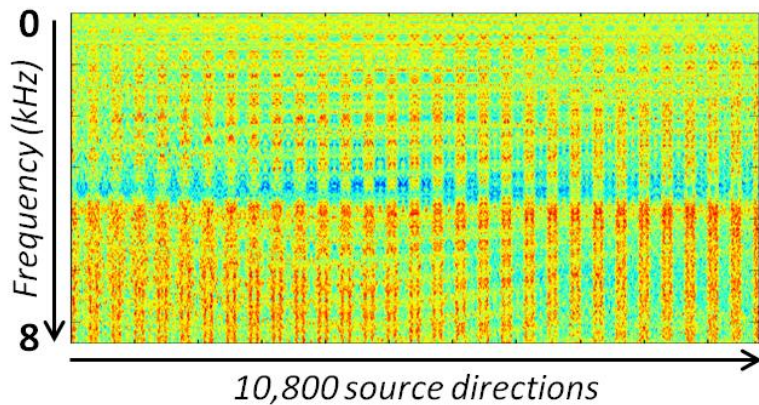
Acoustic Space



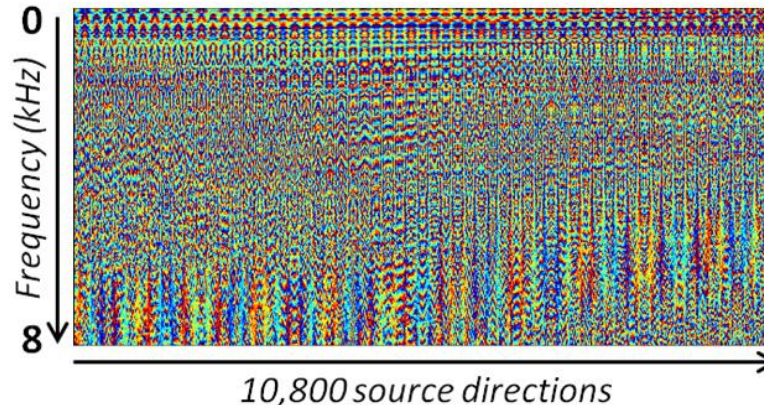
Acoustic Space



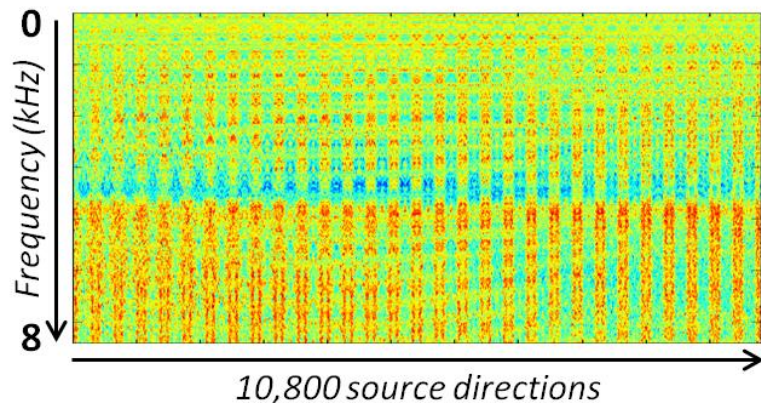
ILD space



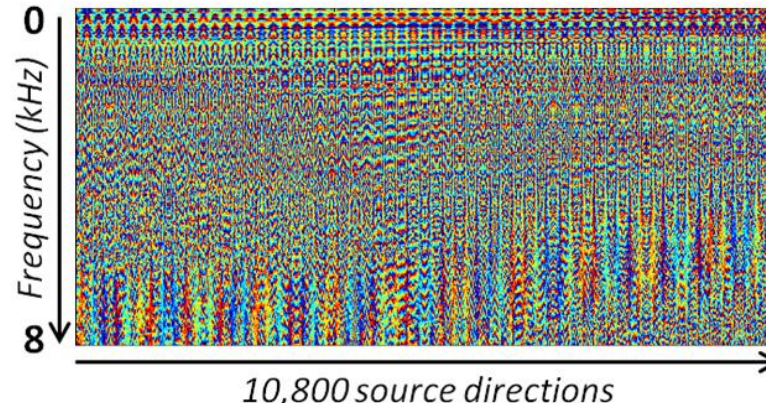
IPD space



ILD space

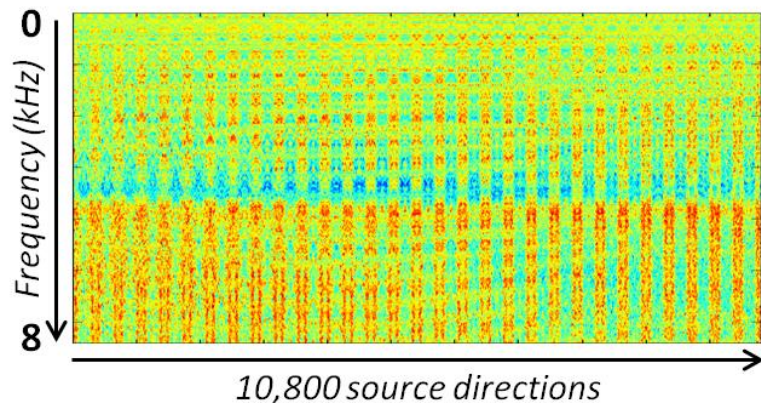


IPD space

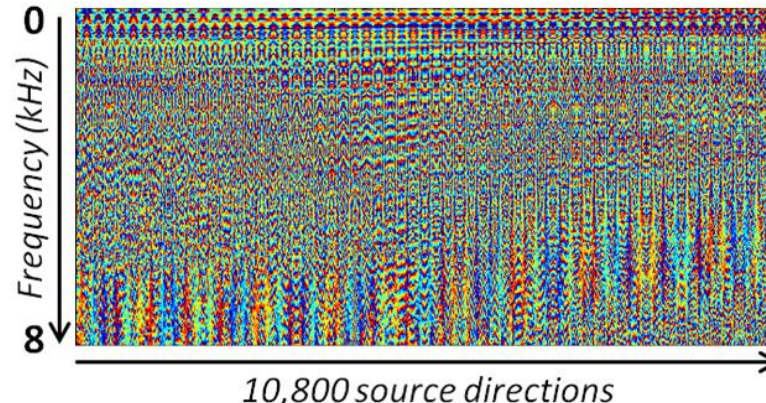


- These high-dimensional representations do not reveal the intrinsic **structure** of acoustic spaces

ILD space

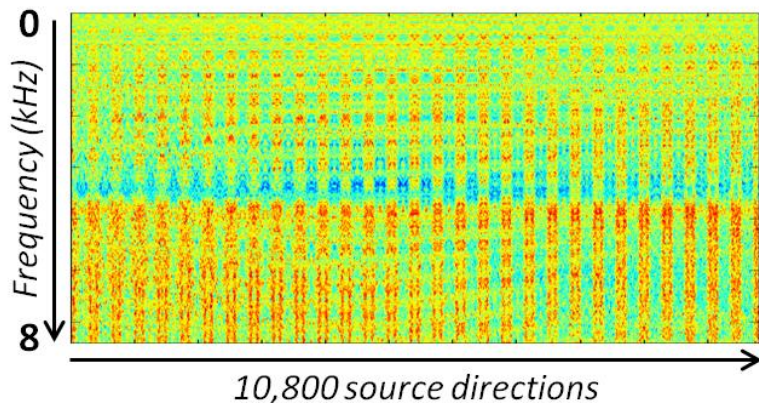


IPD space

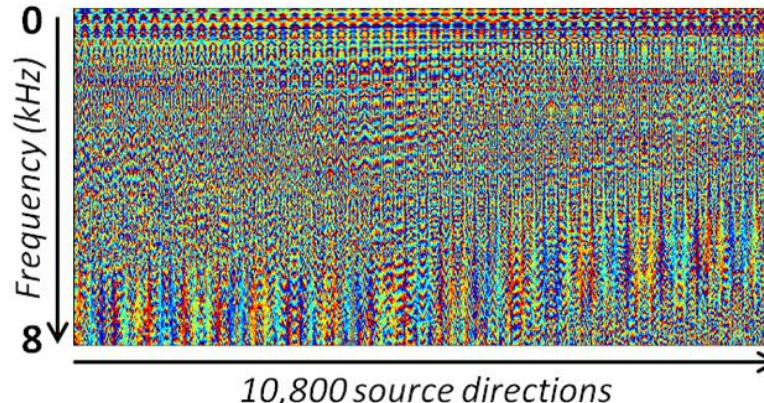


- These high-dimensional representations do not reveal the intrinsic **structure** of acoustic spaces
- We seek a **low-dimensional** representation

ILD space

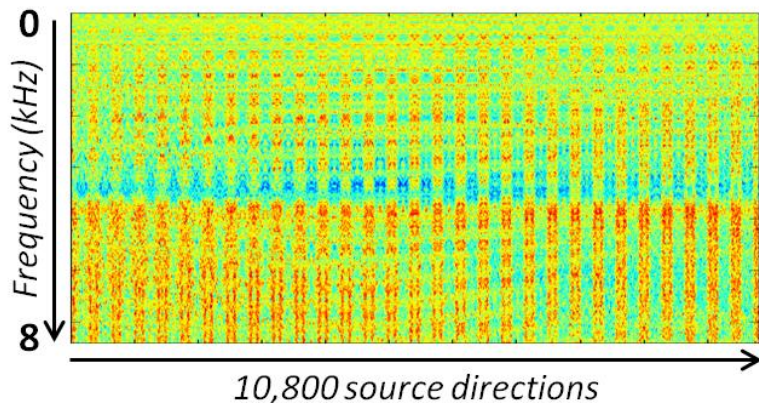


IPD space

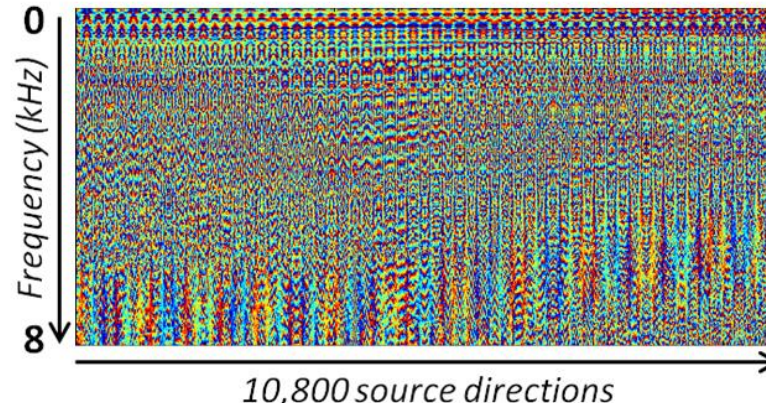


- These high-dimensional representations do not reveal the intrinsic **structure** of acoustic spaces
- We seek a **low-dimensional** representation
- Can be obtained with dimensionality reduction techniques = ***unsupervised learning***

ILD space

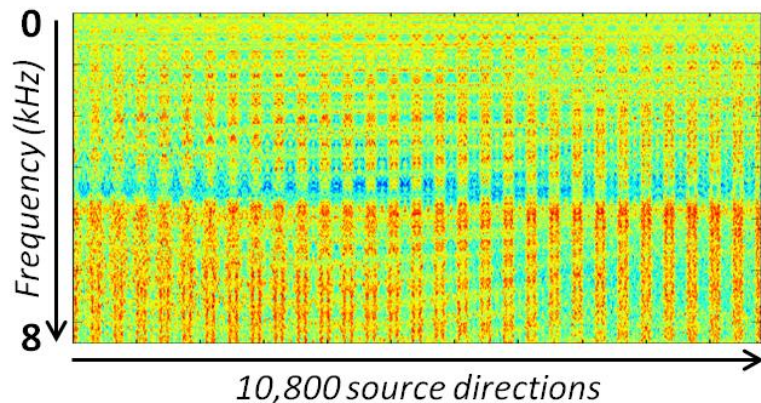


IPD space

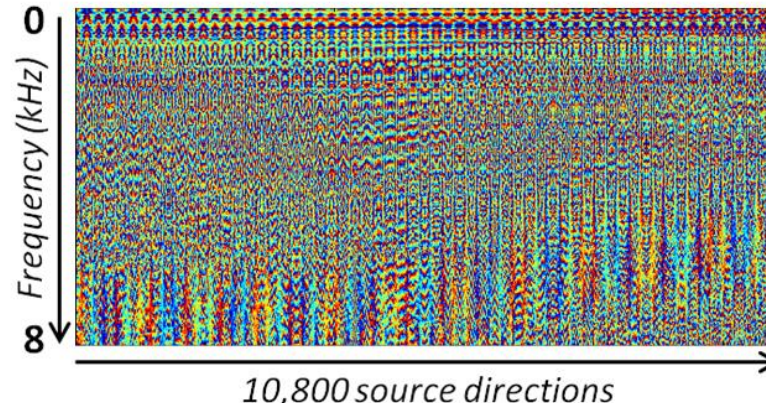


- These high-dimensional representations do not reveal the intrinsic **structure** of acoustic spaces
- We seek a **low-dimensional** representation
- Can be obtained with dimensionality reduction techniques = ***unsupervised learning***
- **Map** data onto a lower dimensional space, using high dimensional data only

ILD space

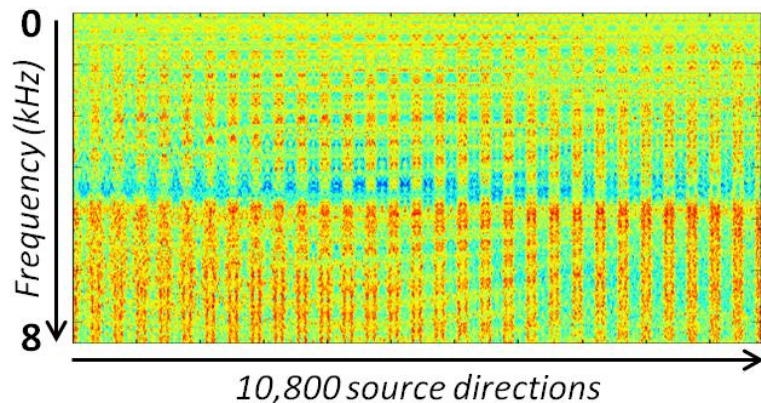


IPD space

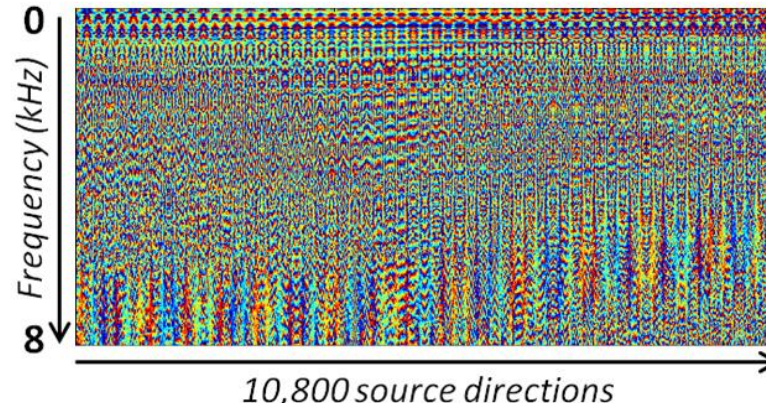


- These high-dimensional representations do not reveal the intrinsic **structure** of acoustic spaces
- We seek a **low-dimensional** representation
- Can be obtained with dimensionality reduction techniques = ***unsupervised learning***
- **Map** data onto a lower dimensional space, using high dimensional data only
- **PCA**=linear, **manifold learning**=non-linear

ILD space



IPD space

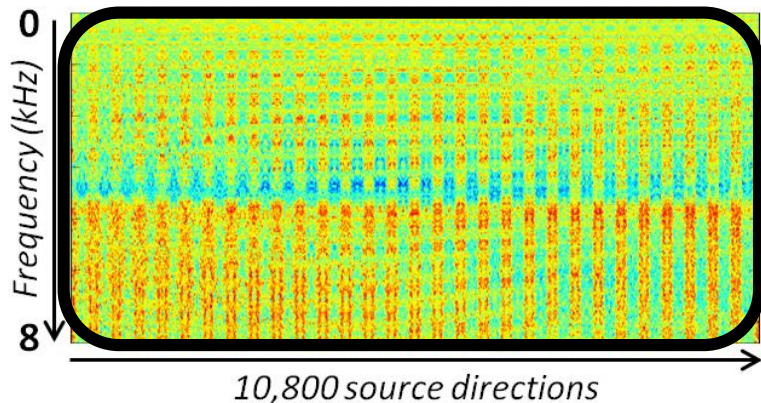


Non-Linear dimensionality reduction (LTSA)

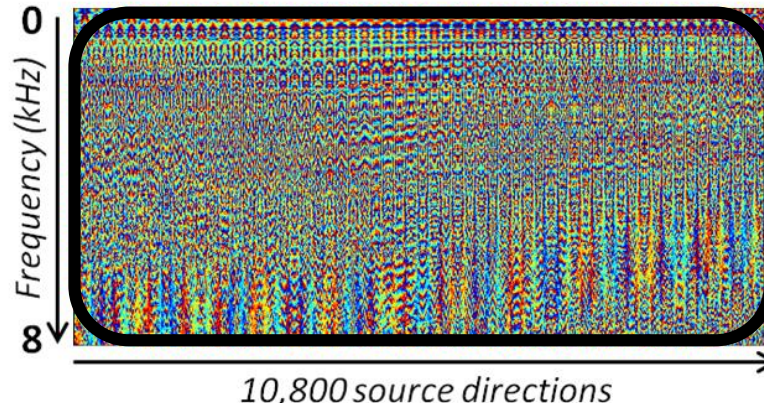
Zhang & Zha (2004)

- Apply PCA locally around each point
- Align globally each representation

ILD space



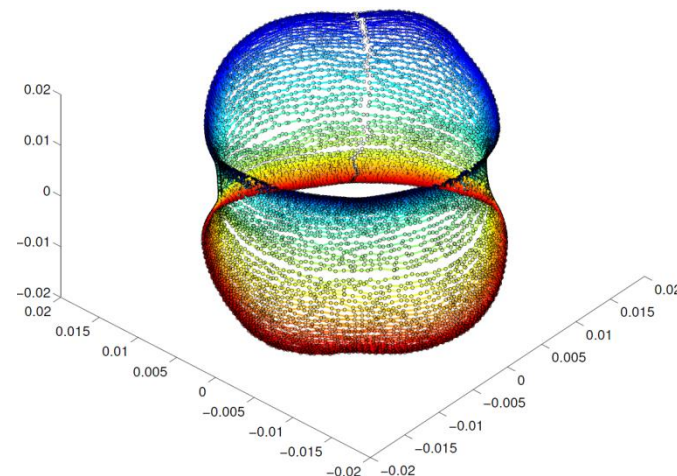
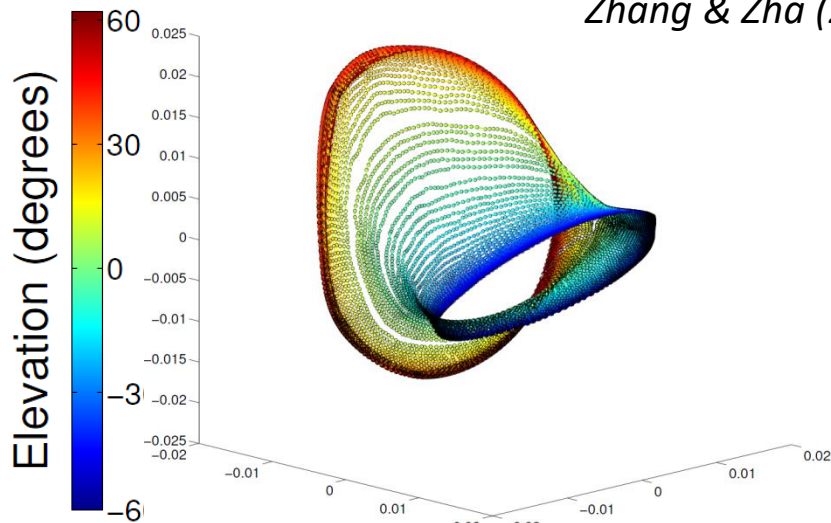
IPD space



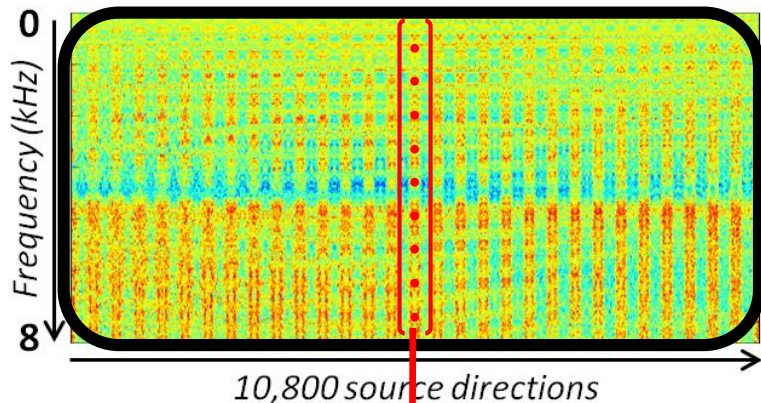
Non-Linear dimensionality reduction (LTSA)

Zhang & Zha (2004)

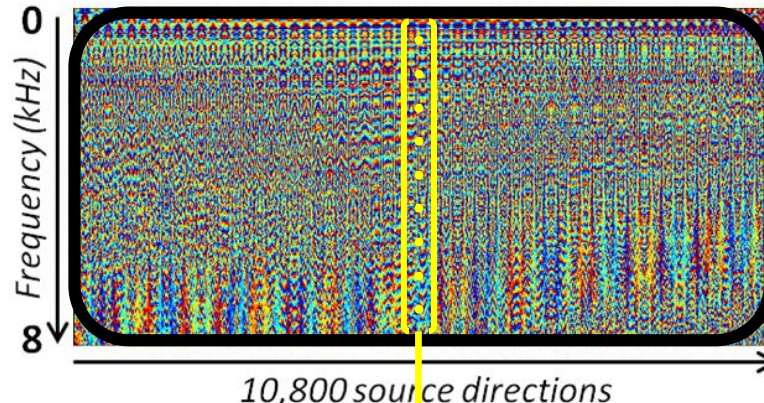
- Apply PCA locally around each point
- Align globally each representation



ILD space



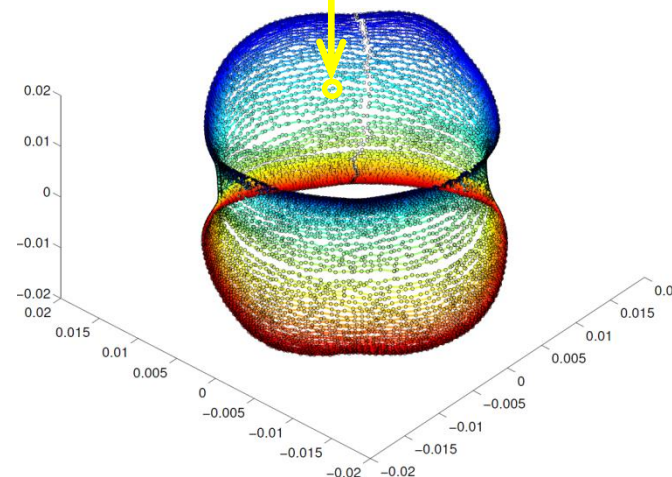
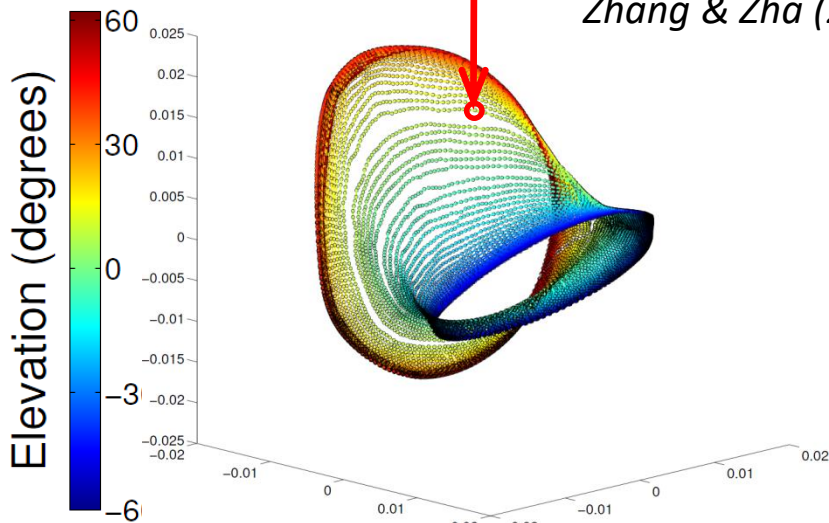
IPD space



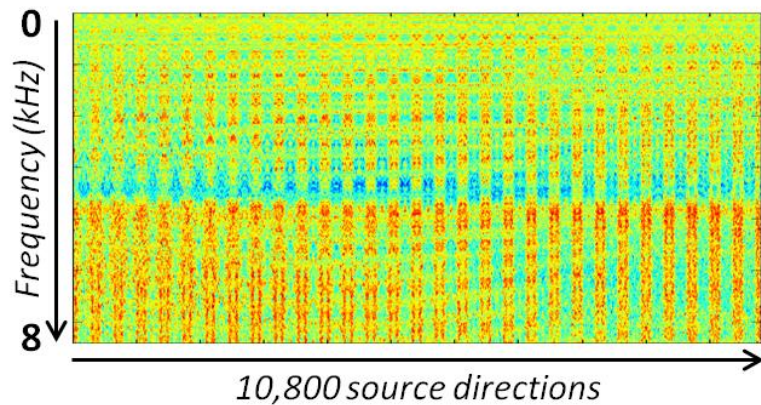
Non-Linear dimensionality reduction (LTSA)

Zhang & Zha (2004)

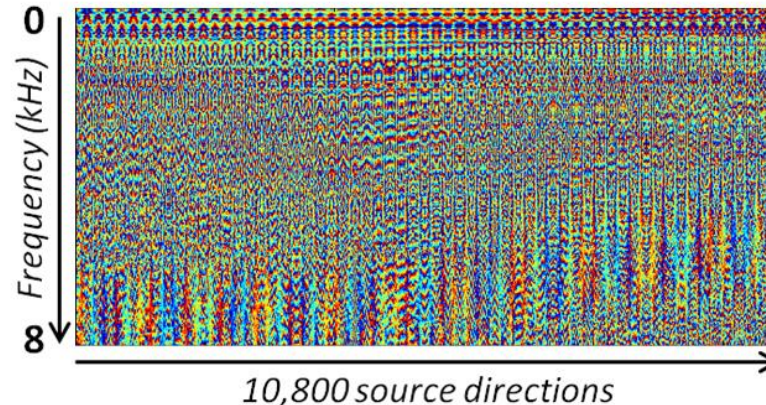
- Apply PCA locally around each point
- Align globally each representation



ILD space

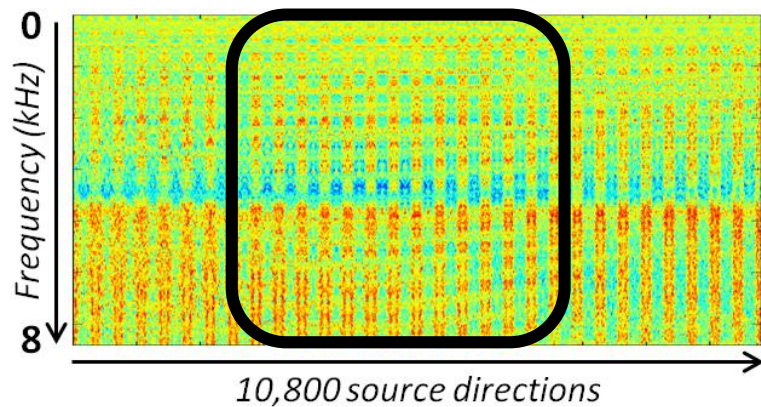


IPD space

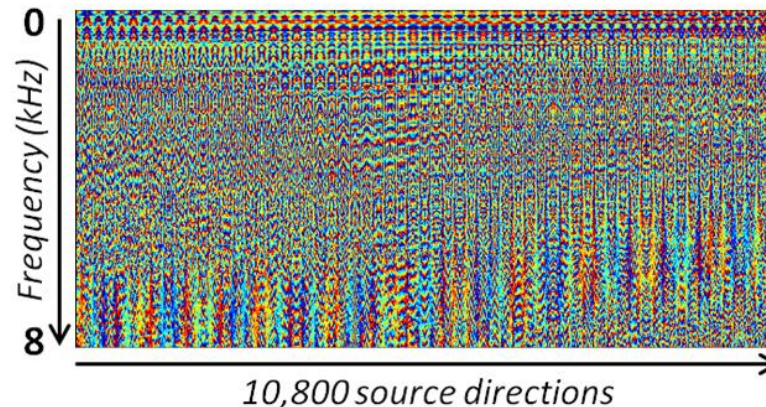


Linear dimensionality reduction (PCA)

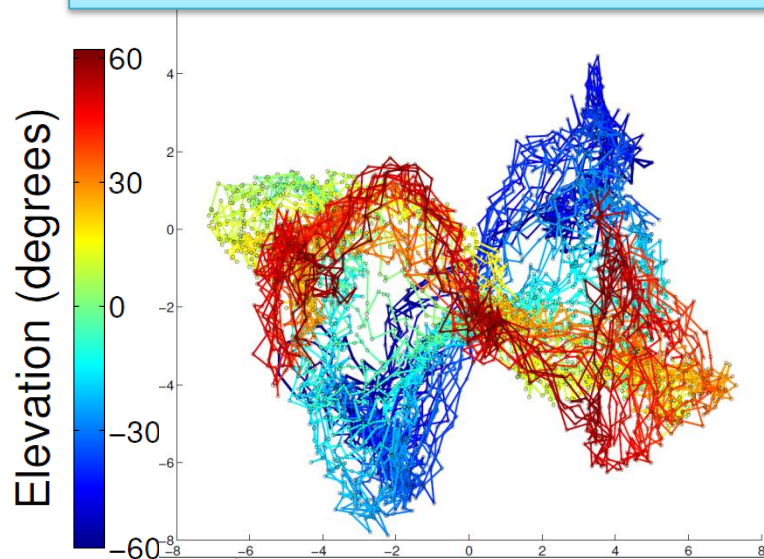
ILD space



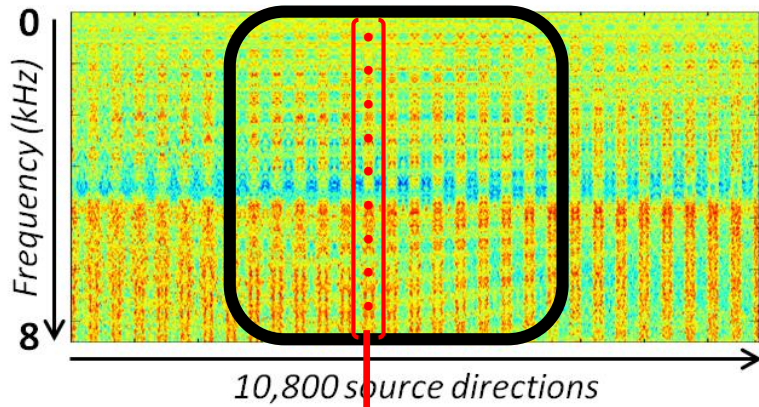
IPD space



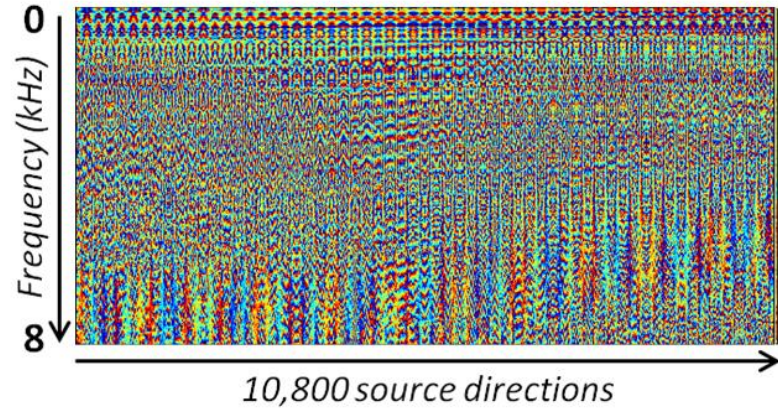
Linear dimensionality reduction (PCA)



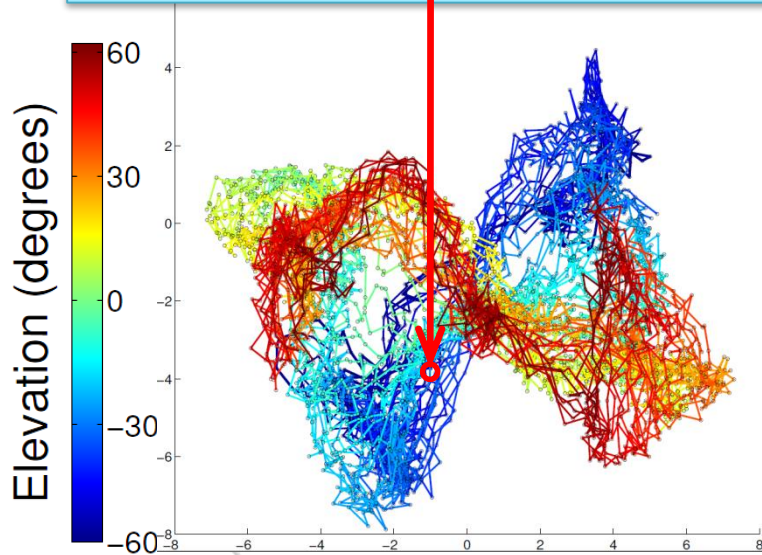
ILD space



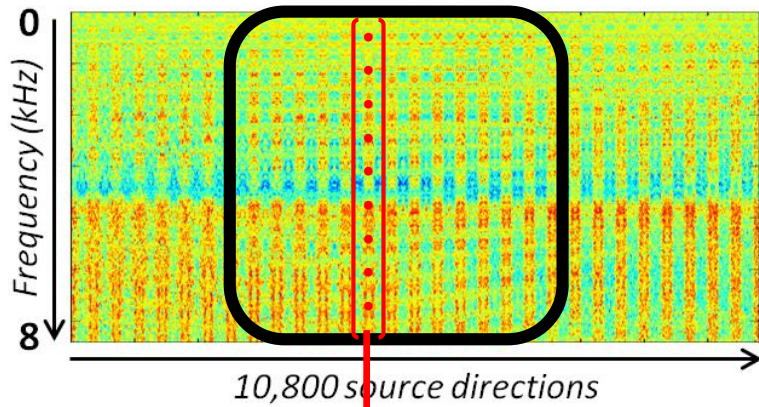
IPD space



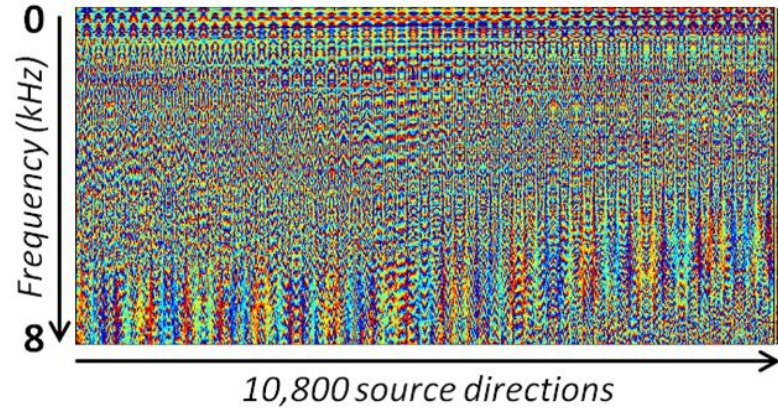
Linear dimensionality reduction (PCA)



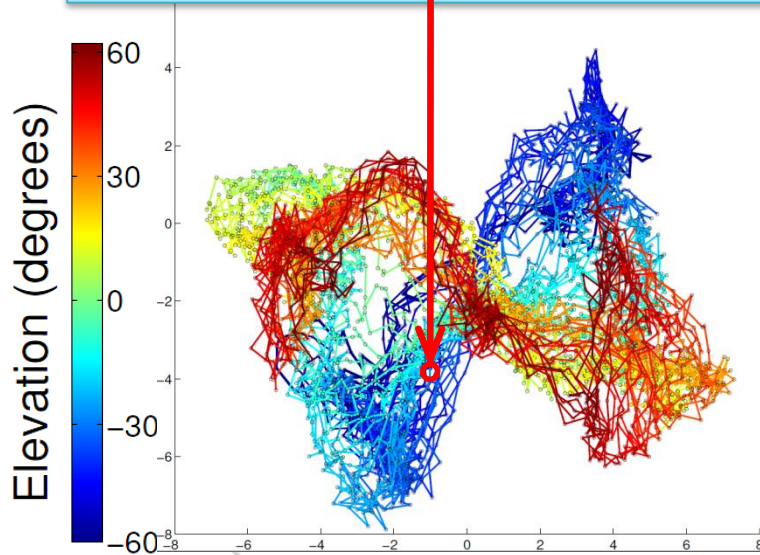
ILD space



IPD space

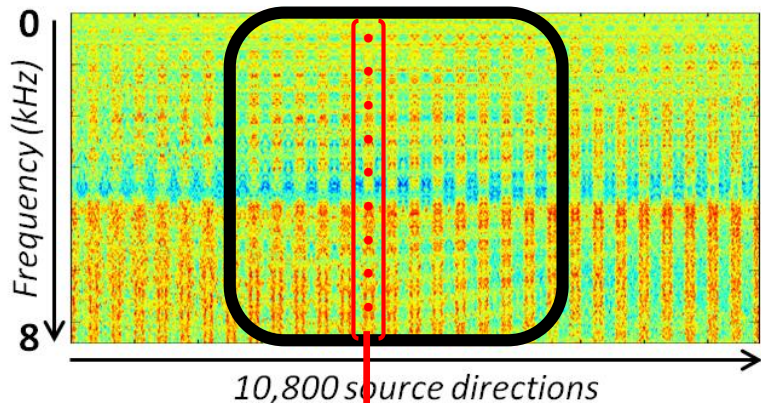


Linear dimensionality reduction (PCA)

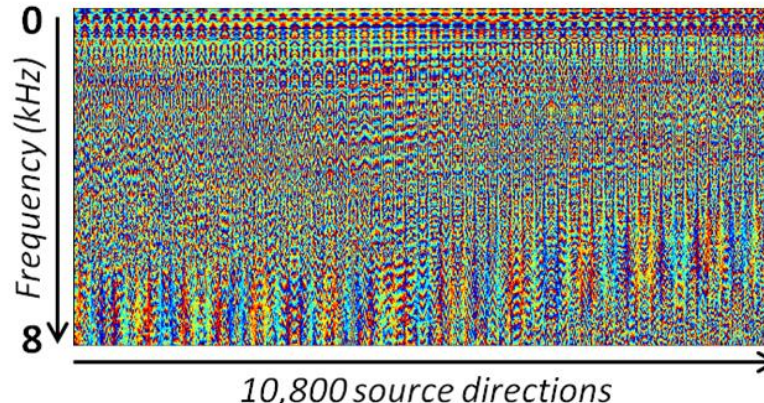


- Very distorted
=> Non-linear structure

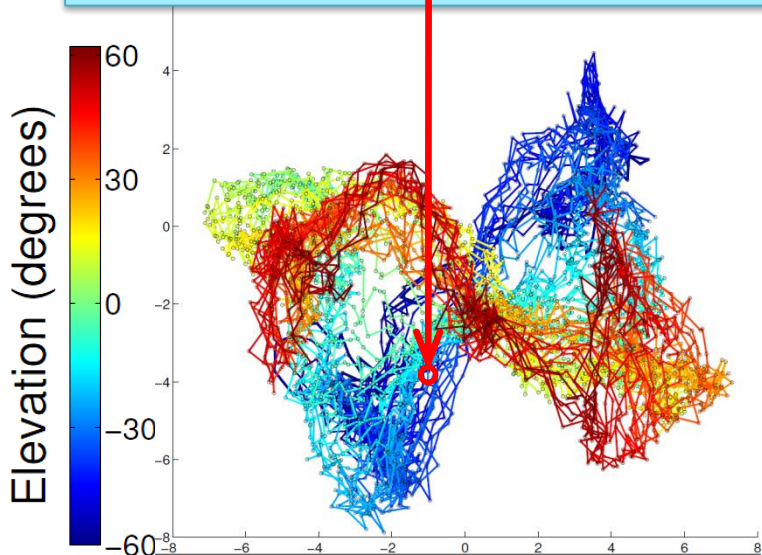
ILD space



IPD space



Linear dimensionality reduction (PCA)

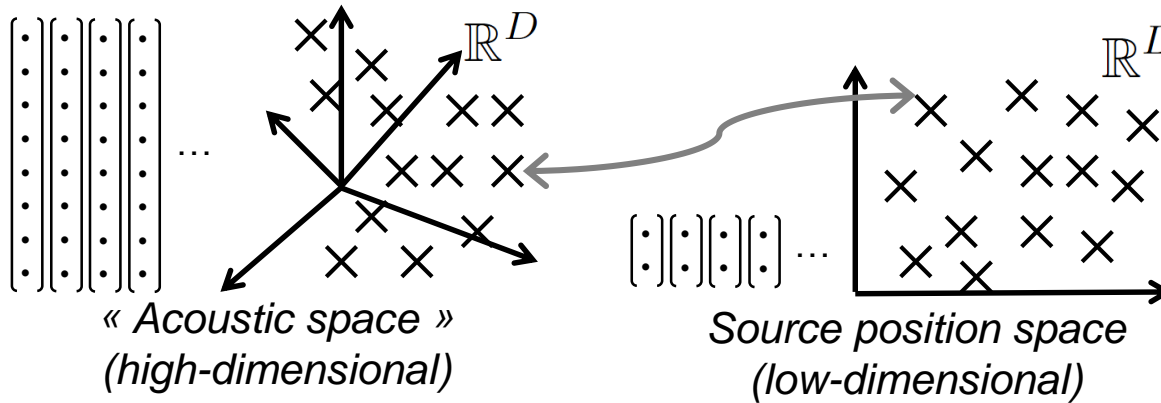


- Very distorted
- => Non-linear structure

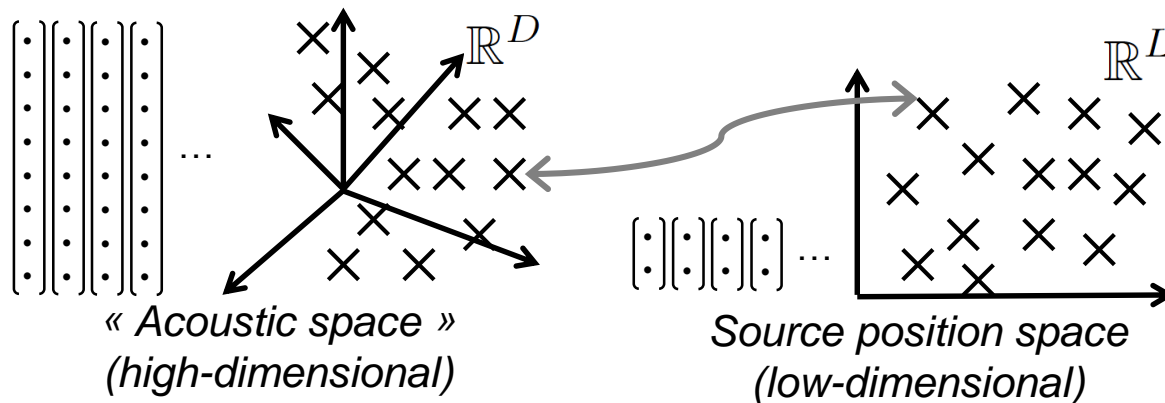
Conclusion on acoustic spaces

- Non-linear but locally-linear
- Lie on a low-dimensional manifold parameterized by source positions

1. Use associated vectors as training data:

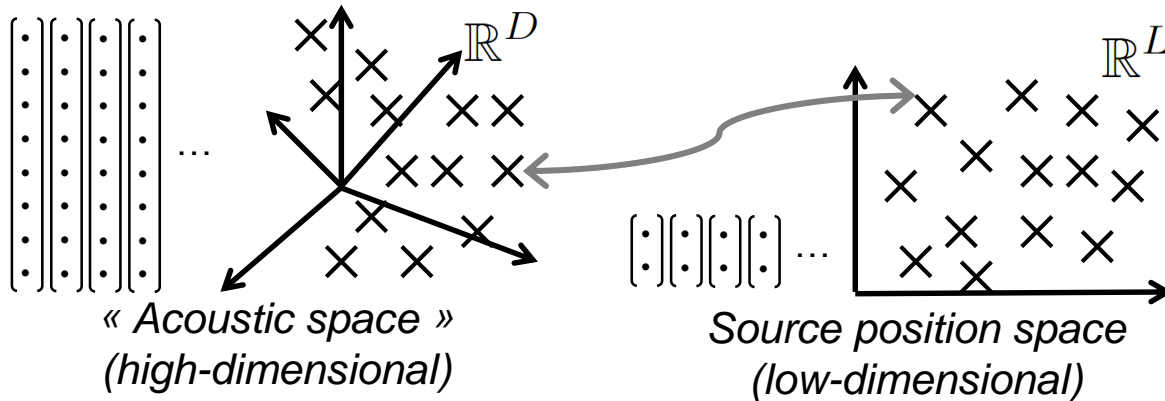


1. Use associated vectors as training data:



How to **collect** training data?

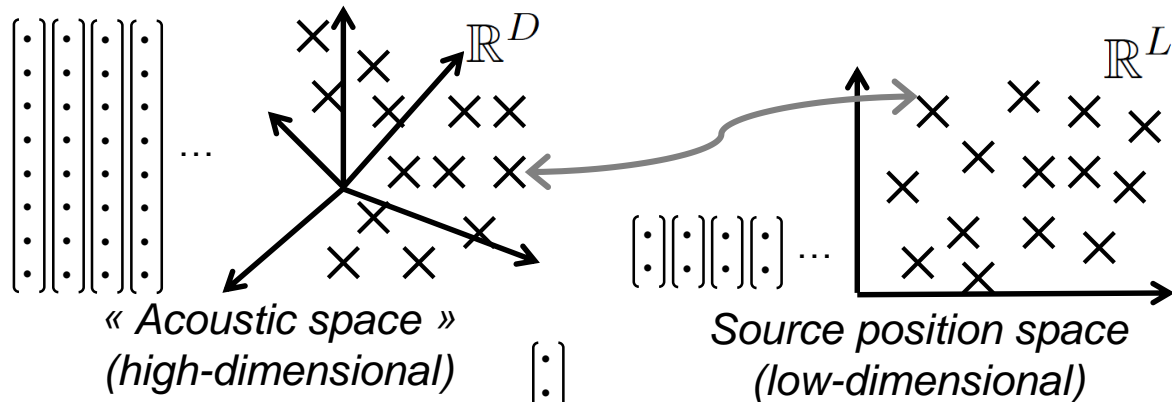
1. Use associated vectors as training data:



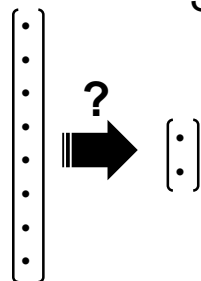
How to **collect** training data?

What is the **structure** of the acoustic space?

1. Use associated vectors as training data:



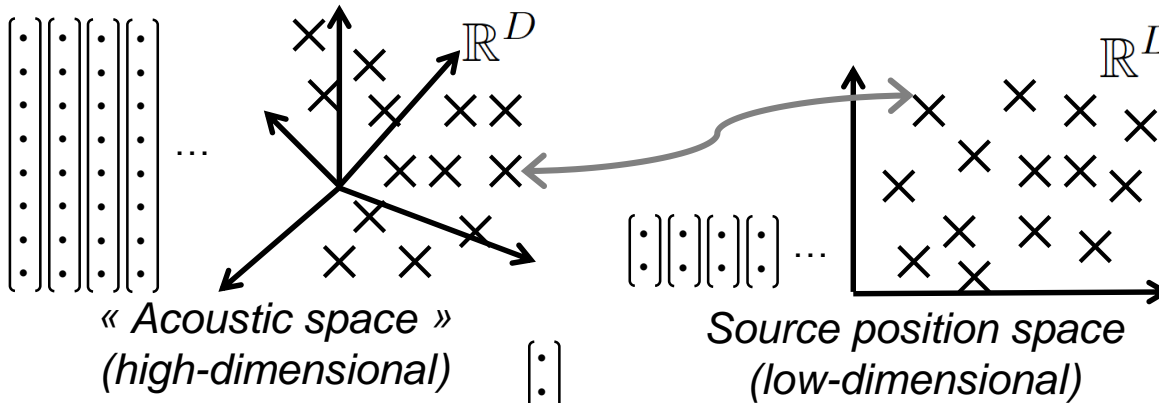
2. Learn the mapping:



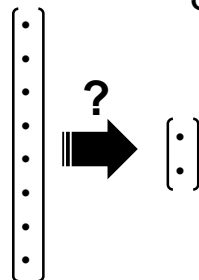
How to **collect** training data?

What is the **structure** of the acoustic space?

1. Use associated vectors as training data:



2. Learn the mapping:

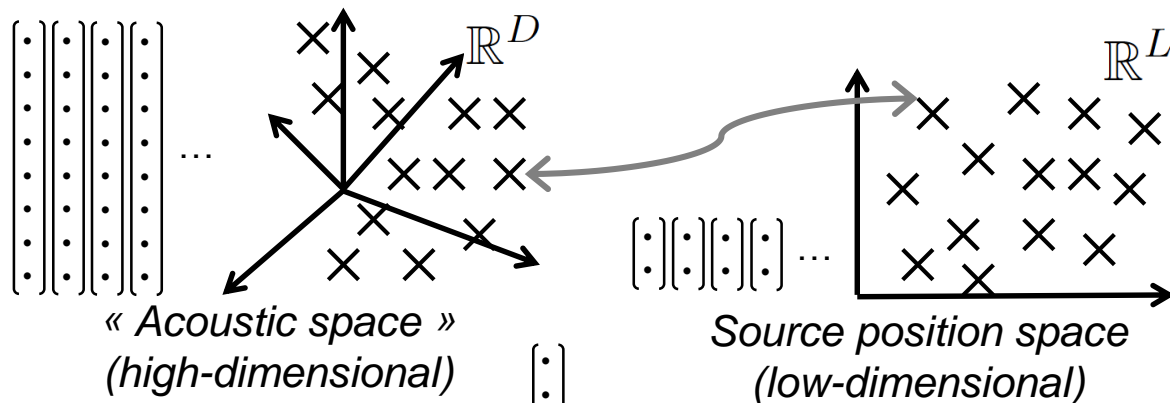


How to **collect** training data?

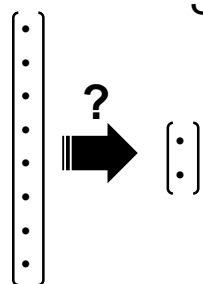
What is the **structure** of the acoustic space?

How to deal with **high-dimensional** input?

1. Use associated vectors as training data:



2. Learn the mapping:



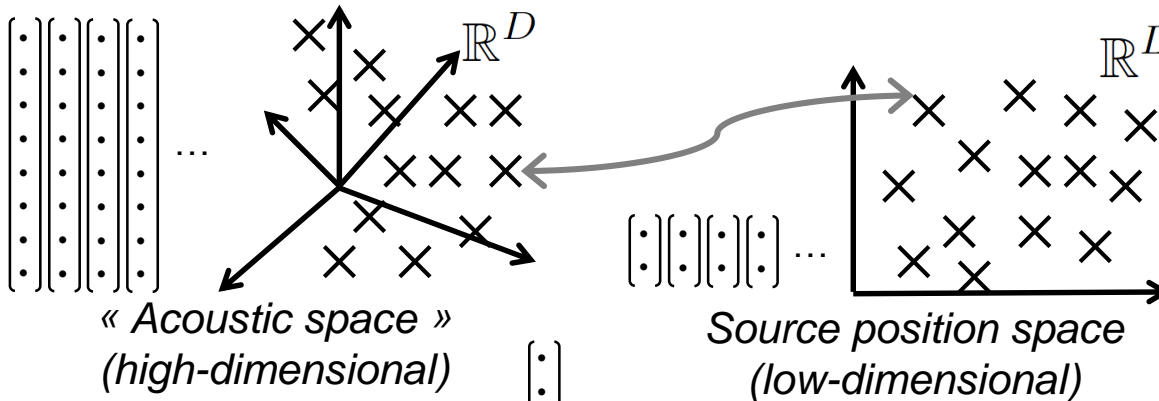
How to **collect** training data?

What is the **structure** of the acoustic space?

How to deal with **high-dimensional** input?

How to estimate a **non-linear** mapping?

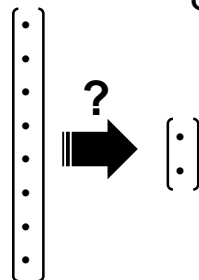
1. Use associated vectors as training data:



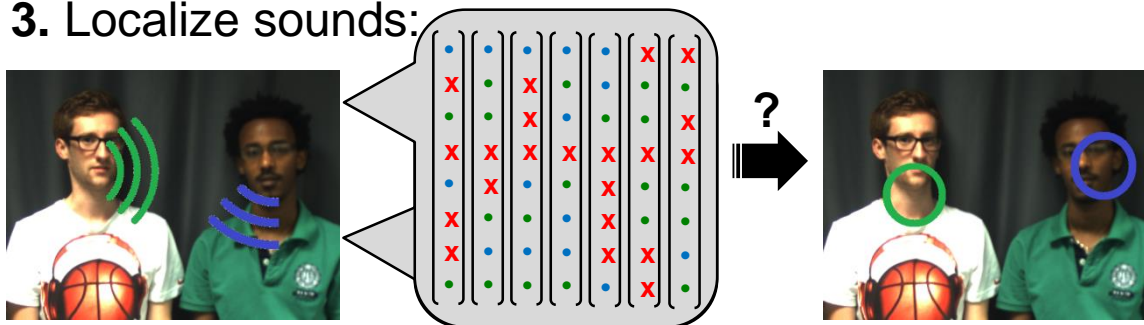
« Acoustic space »
(high-dimensional)

Source position space
(low-dimensional)

2. Learn the mapping:



3. Localize sounds:



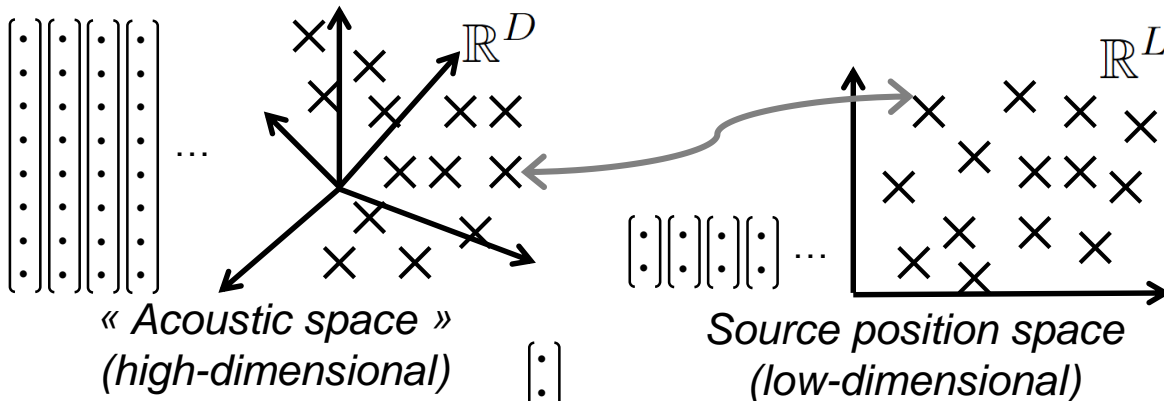
How to **collect**
training data?

What is the **structure** of
the acoustic space?

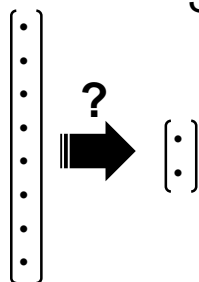
How to deal with **high-**
dimensional input?

How to estimate a **non-**
linear mapping?

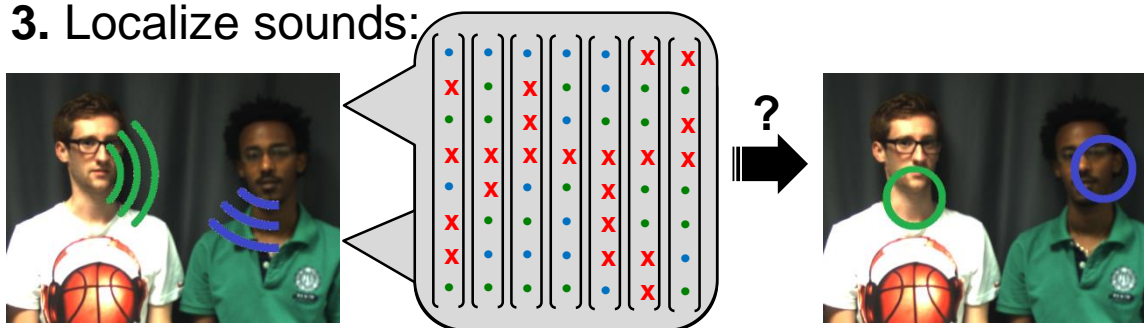
1. Use associated vectors as training data:



2. Learn the mapping:



3. Localize sounds:



How to **collect** training data?

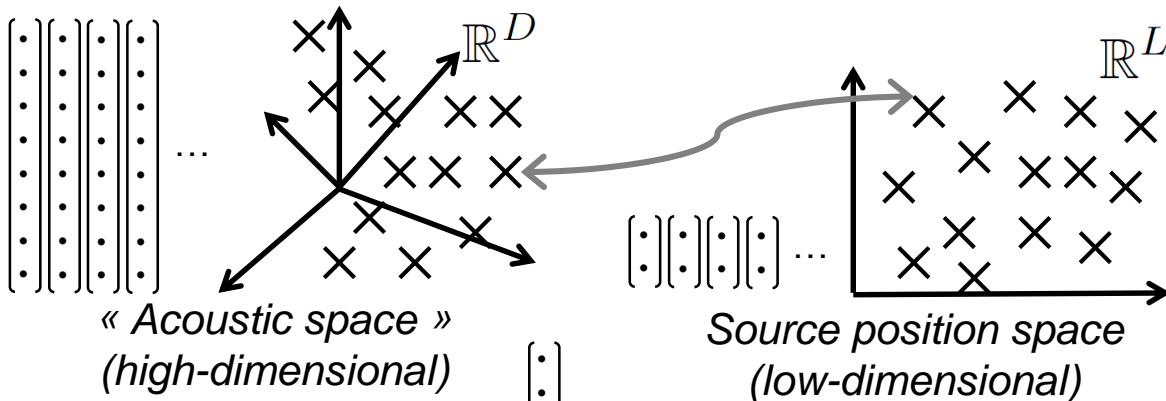
What is the **structure** of the acoustic space?

How to deal with **high-dimensional** input?

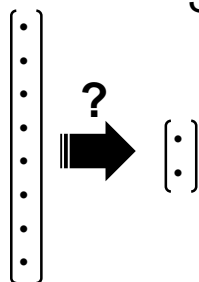
How to estimate a **non-linear** mapping?

How to handle **mixed series** with **missing values** ?

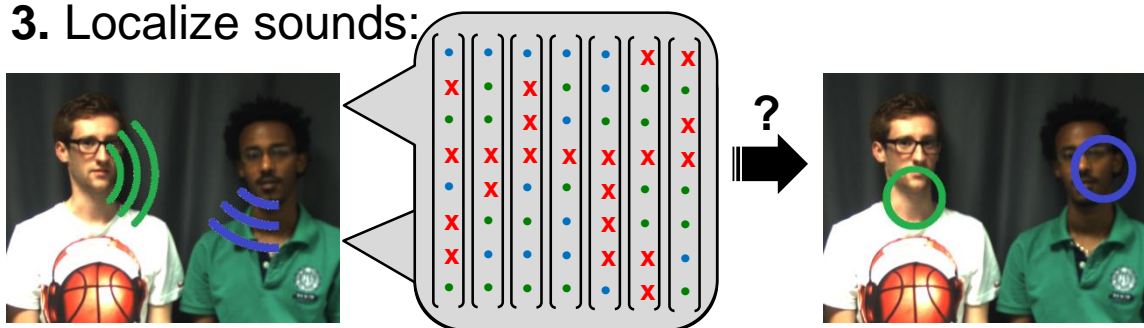
1. Use associated vectors as training data:



2. Learn the mapping:



3. Localize sounds:



How to **collect** training data?

What is the **structure** of the acoustic space?

How to deal with **high-dimensional** input?

How to estimate a **non-linear** mapping?

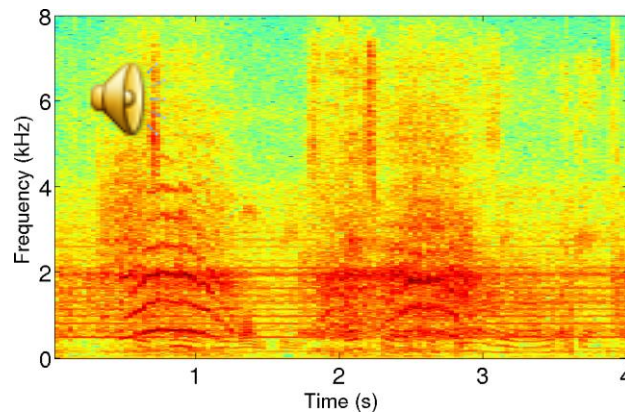
How to handle **mixed series** with **missing values** ?

How to **separate** sound sources?

Or how I kept torturing robots with sounds for a living

Nao Waving

Nao Walking

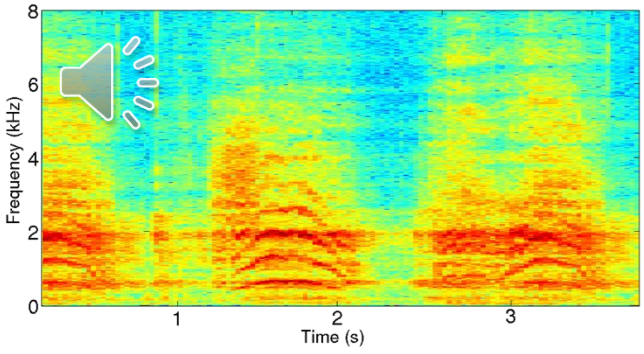


The Post-Doc

Dictionary-based egonoise reduction

≈ 30 seconds of noise only (fan removed with Wiener filtering)

Training Phase

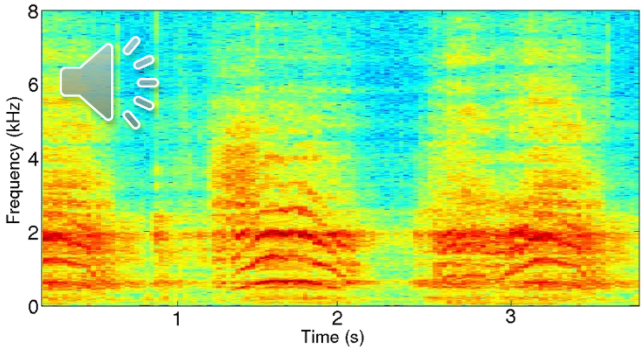


The Post-Doc

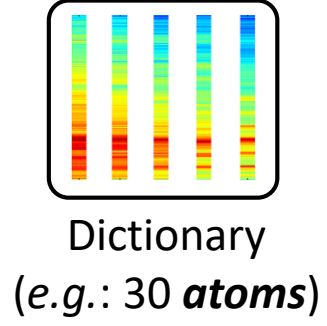
Dictionary-based egonoise reduction

≈ 30 seconds of noise only (fan removed with Wiener filtering)

Training Phase



DL algorithm
(e.g. NMF)

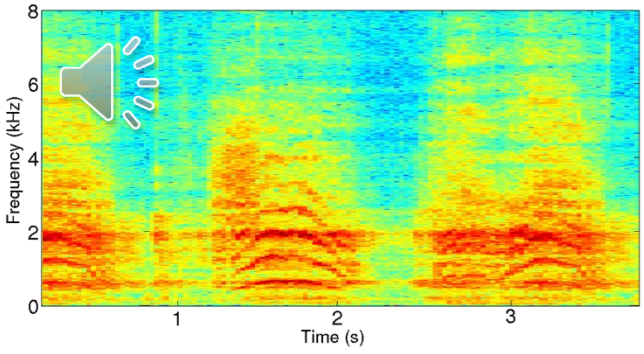


The Post-Doc

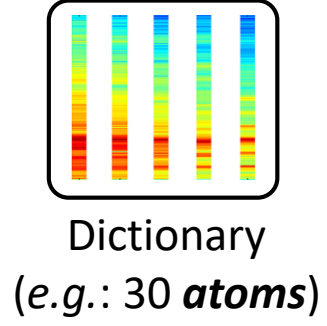
Dictionary-based egonoise reduction

≈ 30 seconds of noise only (fan removed with Wiener filtering)

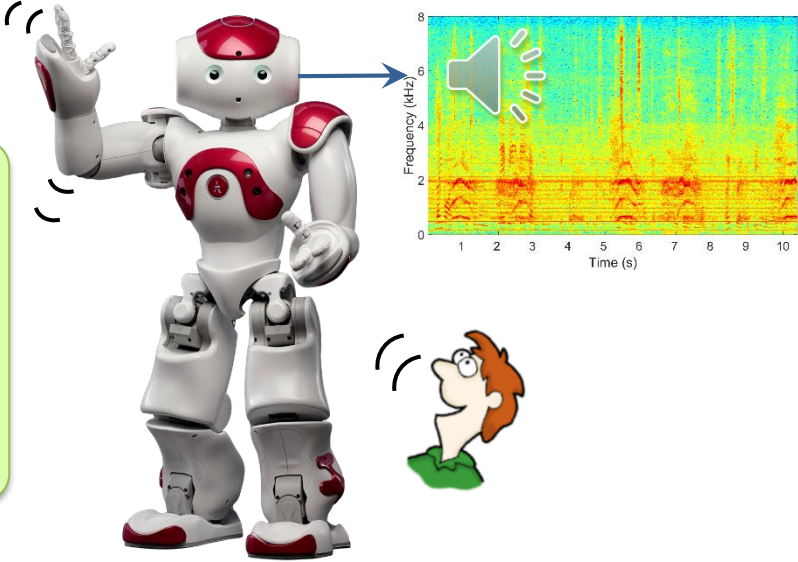
Training Phase



DL algorithm
(e.g. NMF)



Testing Phase

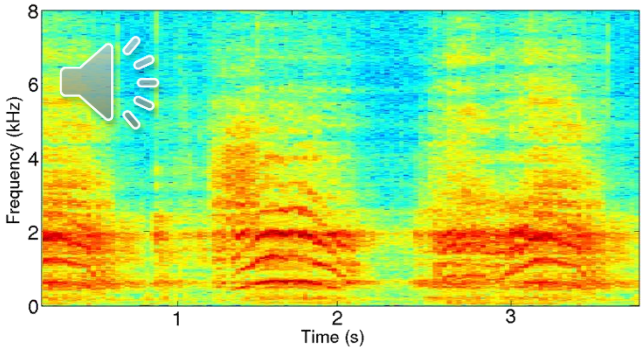


The Post-Doc

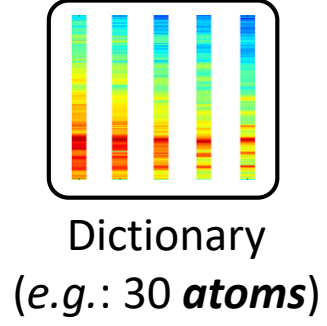
Dictionary-based egonoise reduction

≈ 30 seconds of noise only (fan removed with Wiener filtering)

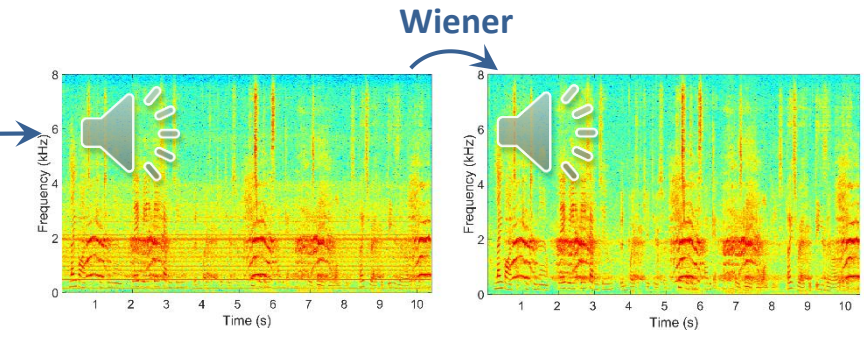
Training Phase



DL algorithm
(e.g. NMF)



Testing Phase

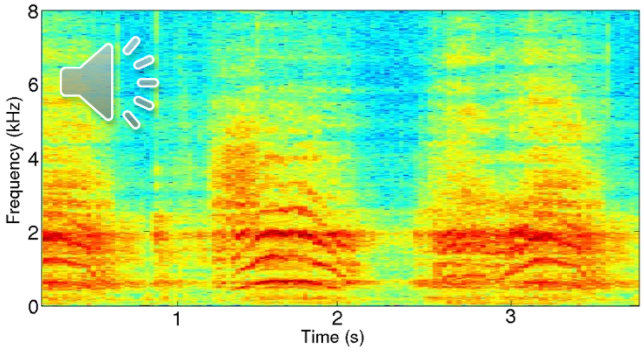


The Post-Doc

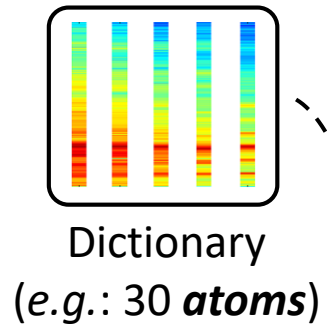
Dictionary-based egonoise reduction

≈ 30 seconds of noise only (fan removed with Wiener filtering)

Training Phase

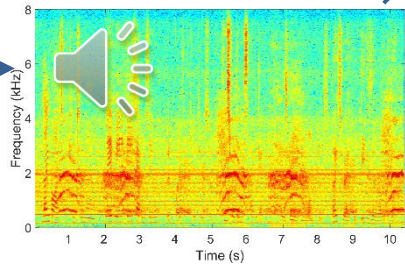


DL algorithm
(e.g. NMF)

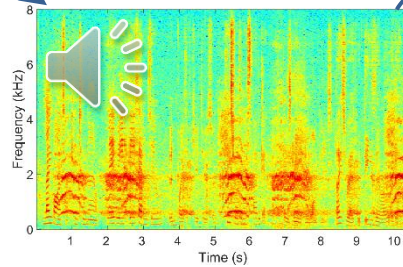


Dictionary
(e.g.: 30 *atoms*)

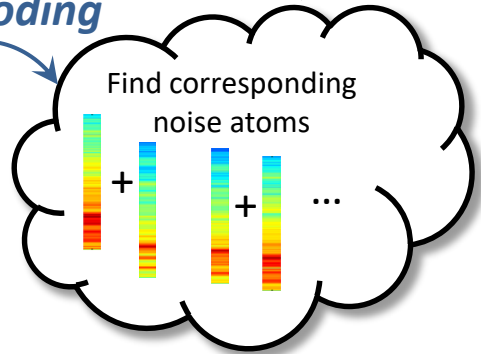
Testing Phase



Wiener



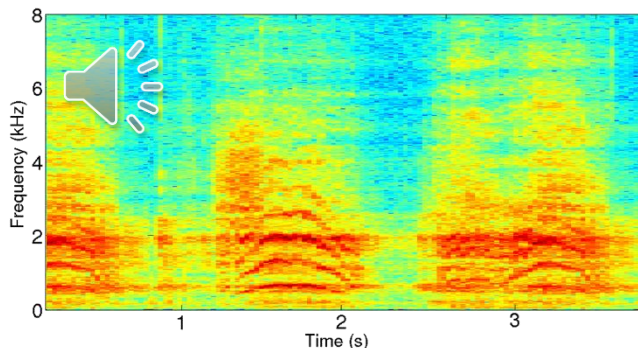
Sparse Coding



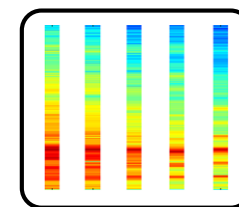
Dictionary-based egonoise reduction

≈ 30 seconds of noise only (fan removed with Wiener filtering)

Training Phase

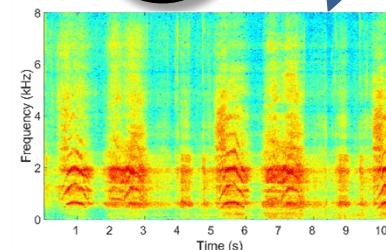
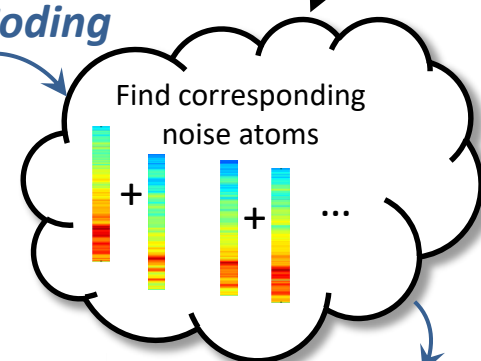
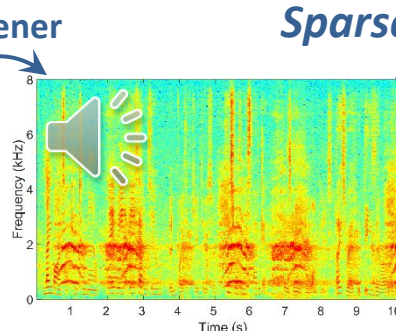
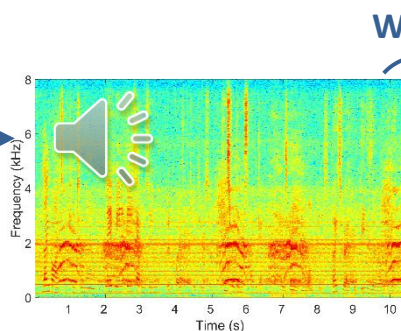
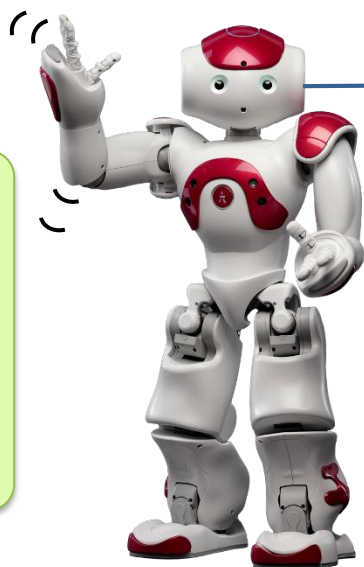


DL algorithm
(e.g. NMF)



Dictionary
(e.g.: 30 *atoms*)

Testing Phase

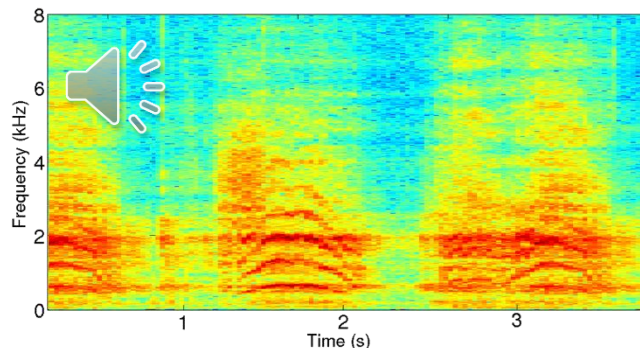


Reconstructed noise

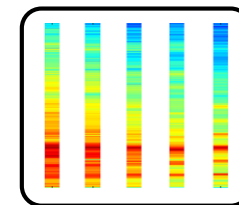
Dictionary-based egonoise reduction

≈ 30 seconds of noise only (fan removed with Wiener filtering)

Training Phase

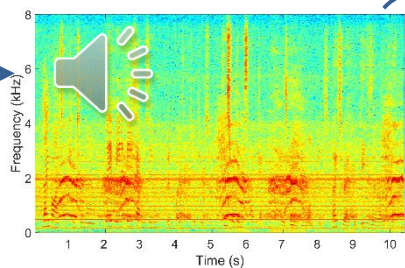


DL algorithm
(e.g. NMF)

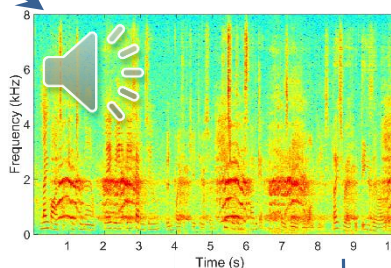


Dictionary
(e.g.: 30 *atoms*)

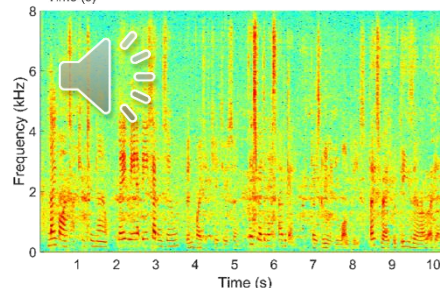
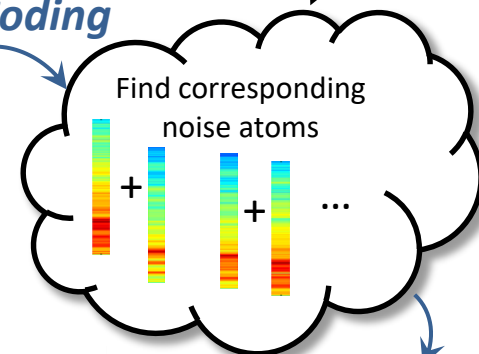
Testing Phase



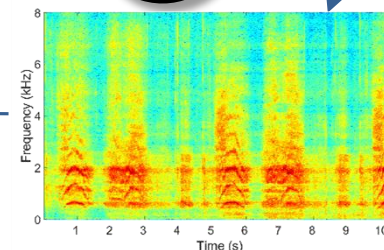
Wiener



Sparse Coding



Cleaned Signal



Reconstructed noise

Ok, let's see on what robot I could release my psychopatic urges this time...

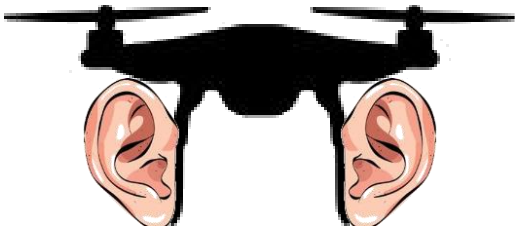
Ok, let's see on what robot I could release my psychopatic urges this time...

« *Oh, hey, I heard you were working with **Drones** ?* »

Ok, let's see on what robot I could release my psychopatic urges this time...

« *Oh, hey, I heard you were working with **Drones** ?* »

« *What do you think of a drone with ears?* »



Ok, let's see on what robot I could release my psychopatic urges this time...

« *Oh, hey, I heard you were working with **Drones** ?* »

« *What do you think of a drone with ears?* »



The DREGON dataset: Drone Audition for Search & Rescue

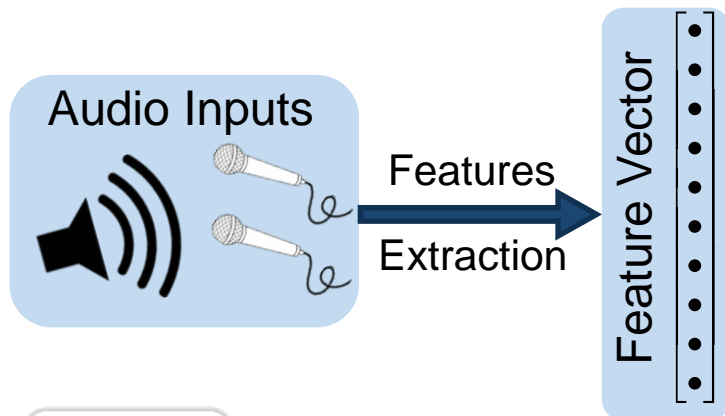
Or the sudden realization that gathering audio data to train robots is a massive pain in the... is very impractical.

Or the sudden realization that gathering audio data to train robots is a massive pain in the... is very impractical.

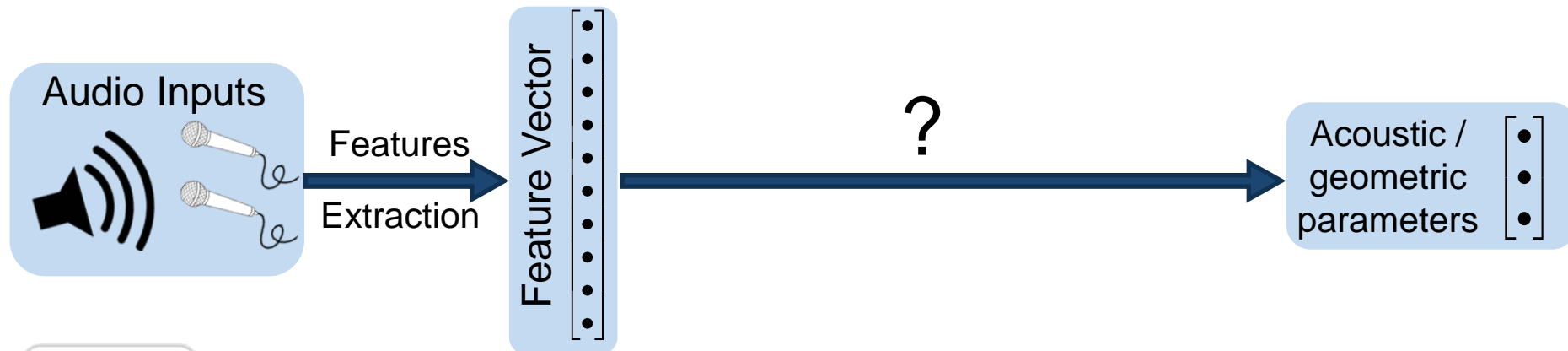
Audio Inputs



Or the sudden realization that gathering audio data to train robots is a massive pain in the... is very impractical.

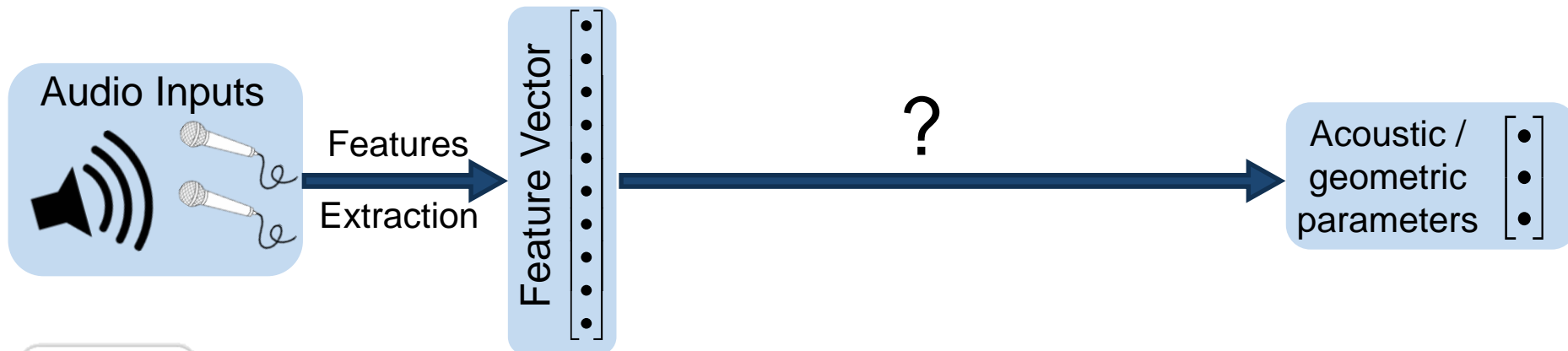


Or the sudden realization that gathering audio data to train robots is a massive pain in the... is very impractical.

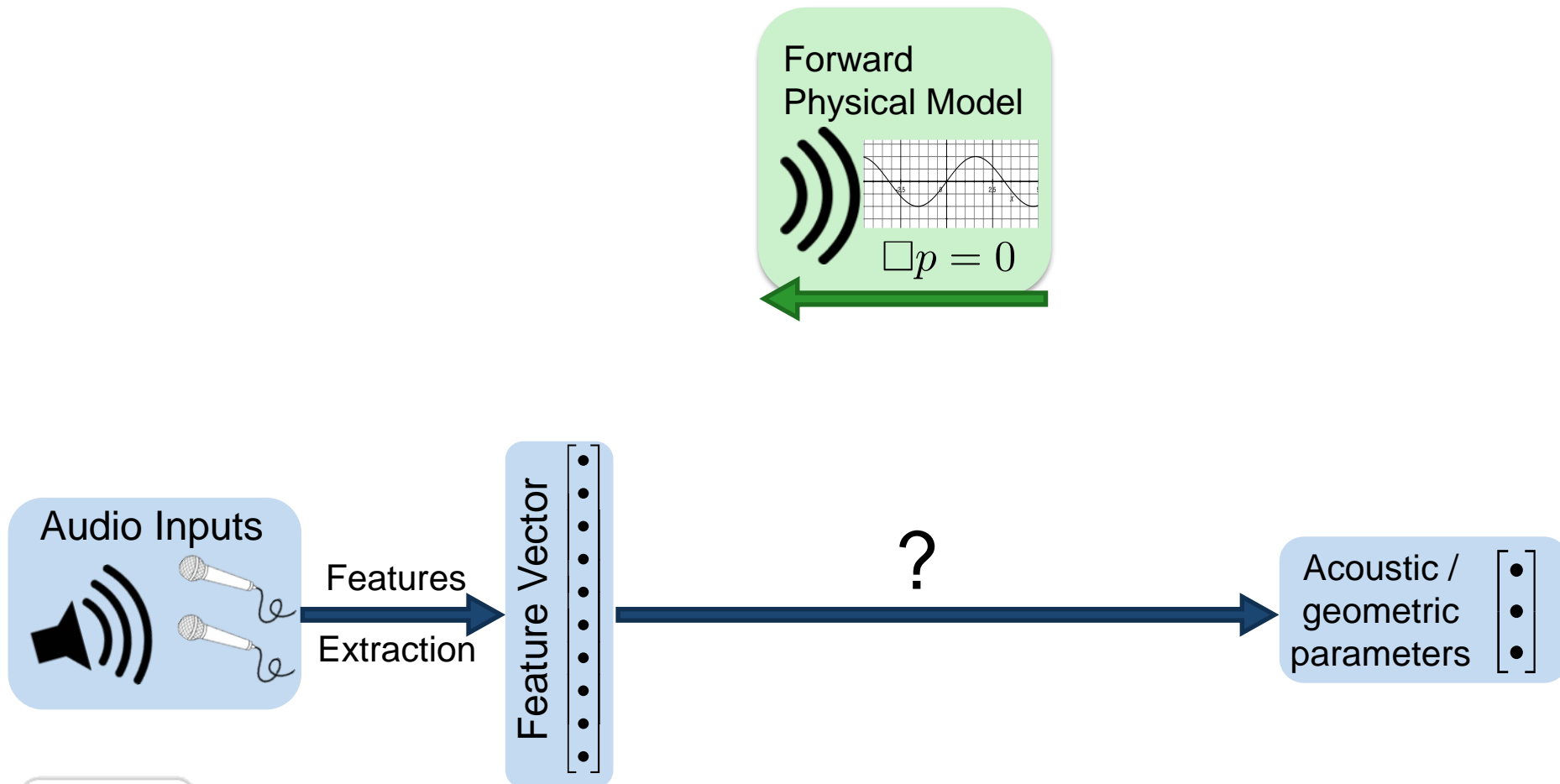


Present time, Inria Nancy

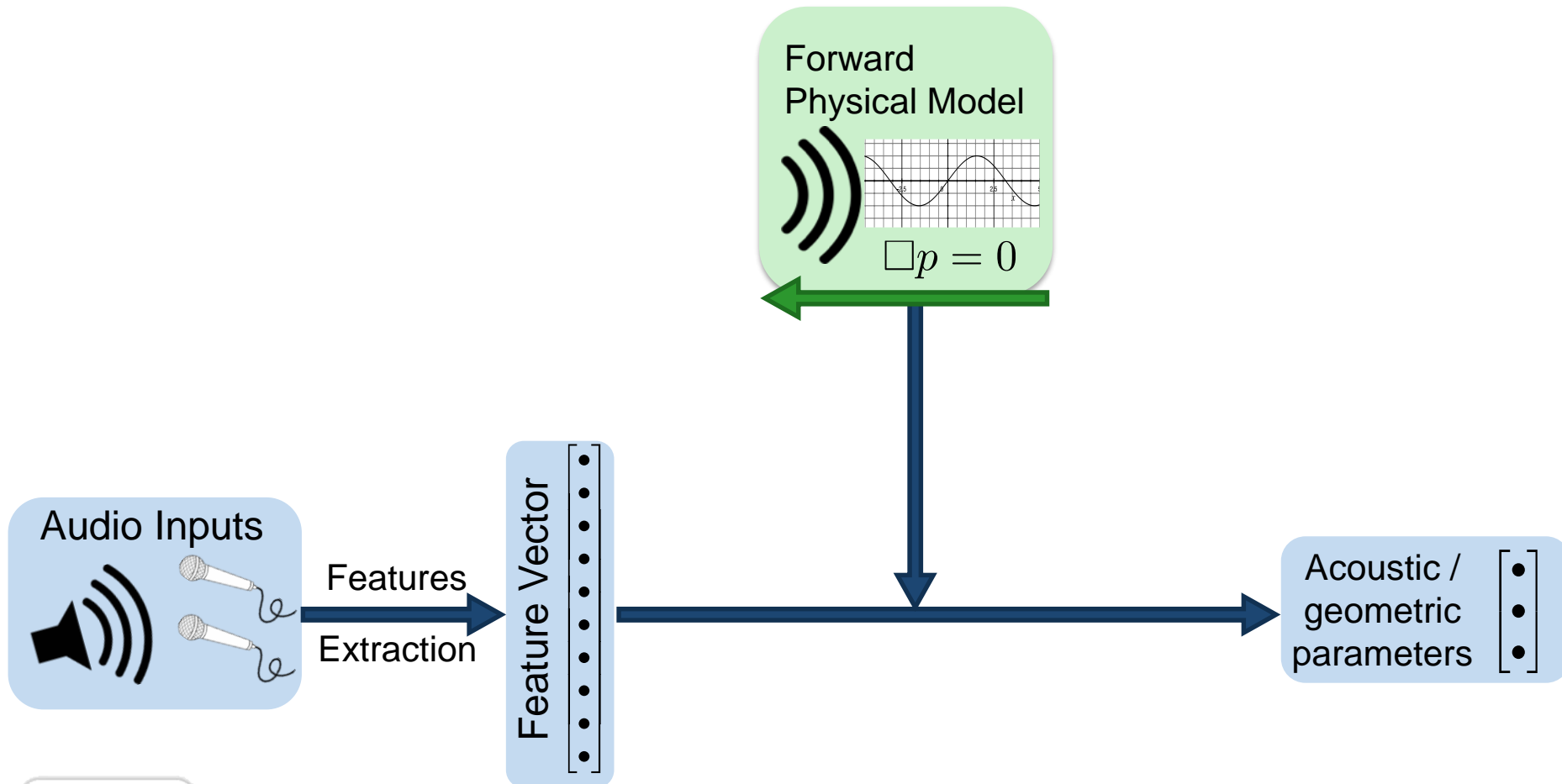
a) Physics-Driven / Traditional Approaches



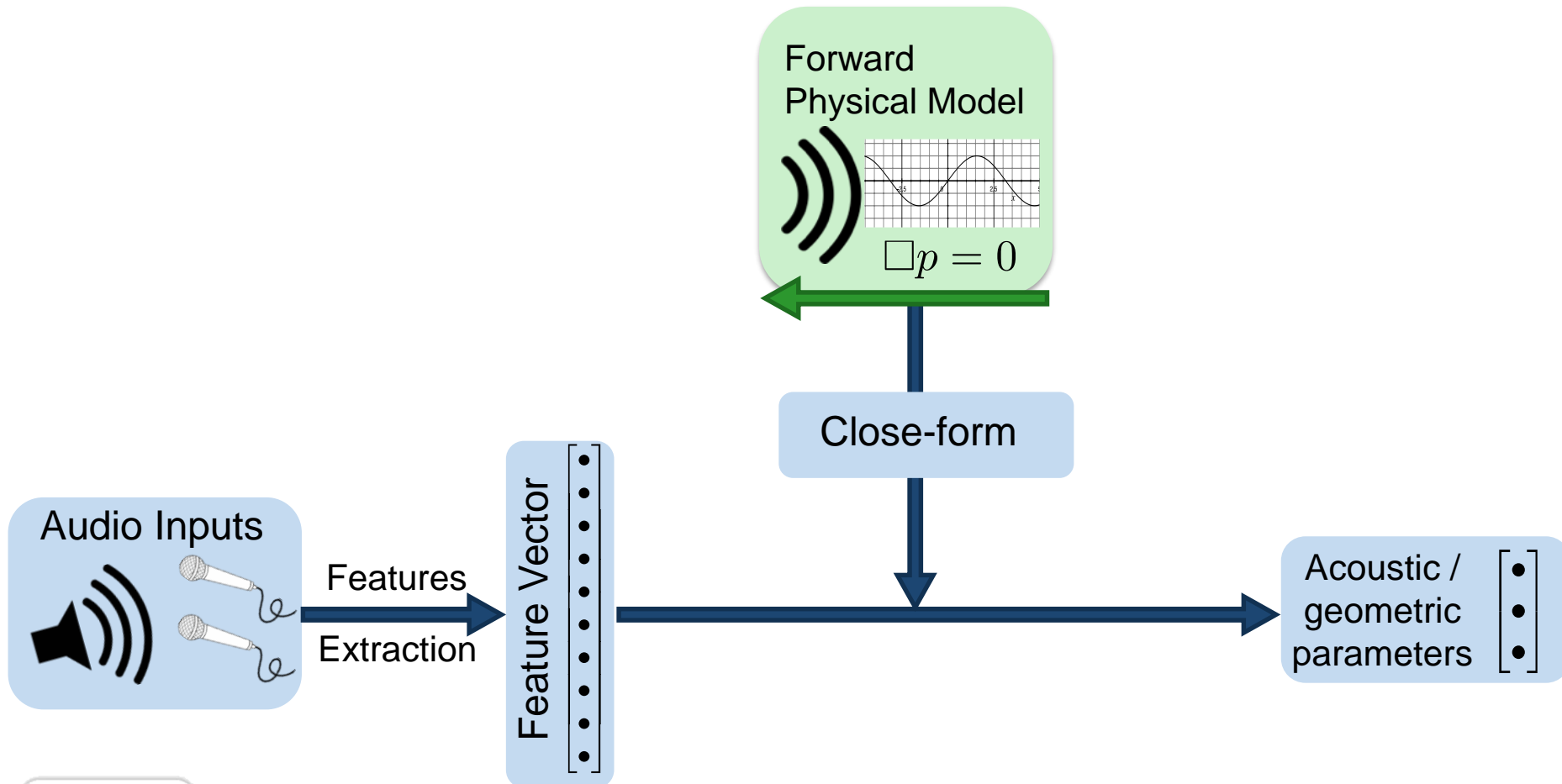
a) Physics-Driven / Traditional Approaches



a) Physics-Driven / Traditional Approaches

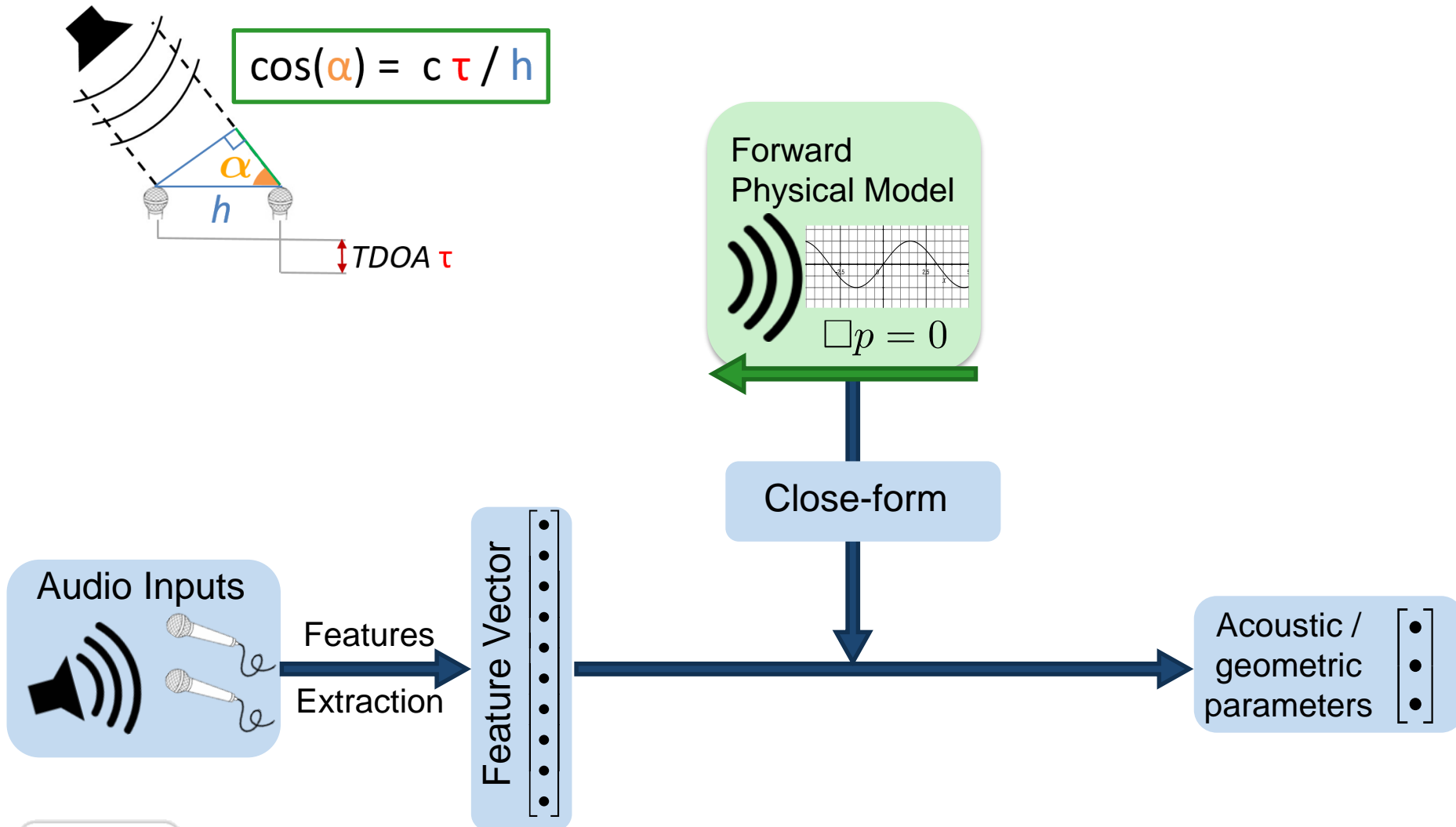


a) Physics-Driven / Traditional Approaches



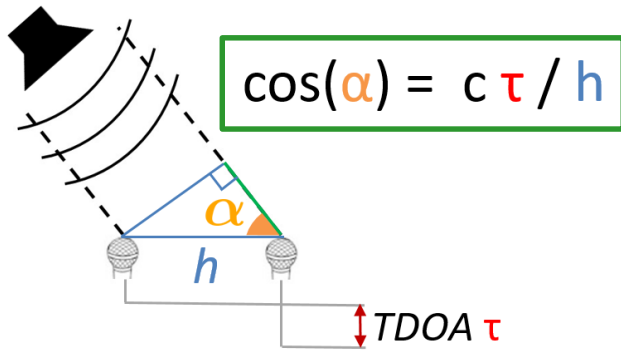
Present time, Inria Nancy

a) Physics-Driven / Traditional Approaches



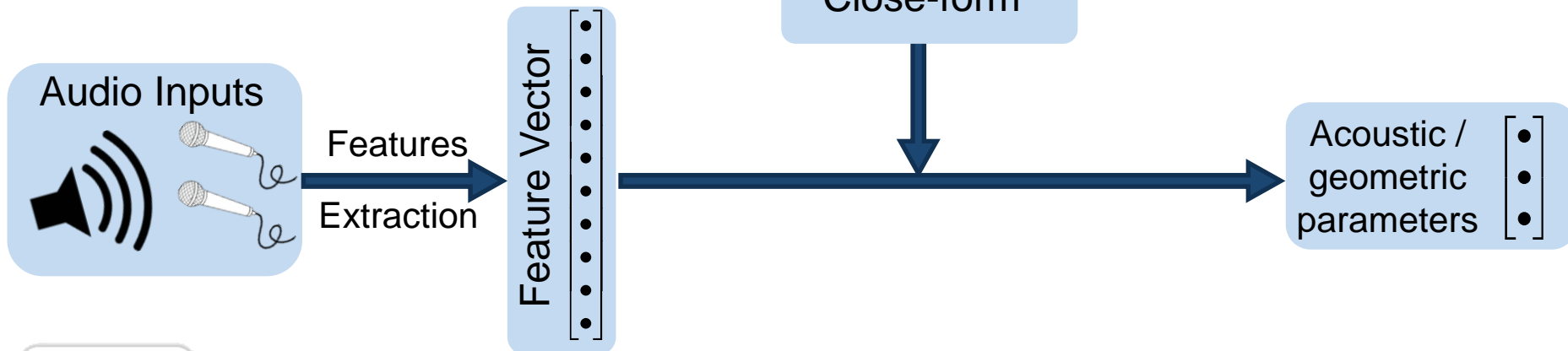
Present time, Inria Nancy

a) Physics-Driven / Traditional Approaches



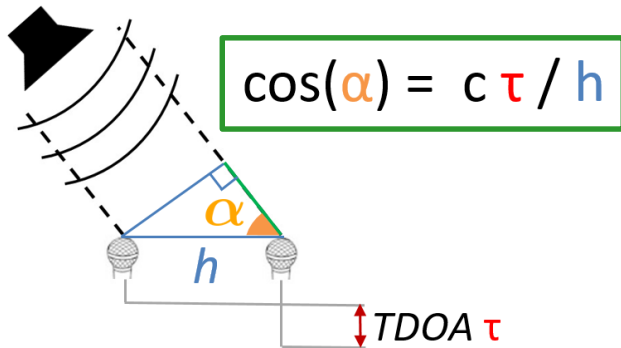
Sabine's law:

$$RT_{60}(b) \approx 0.16 \frac{V}{S \bar{\alpha}(b)}$$



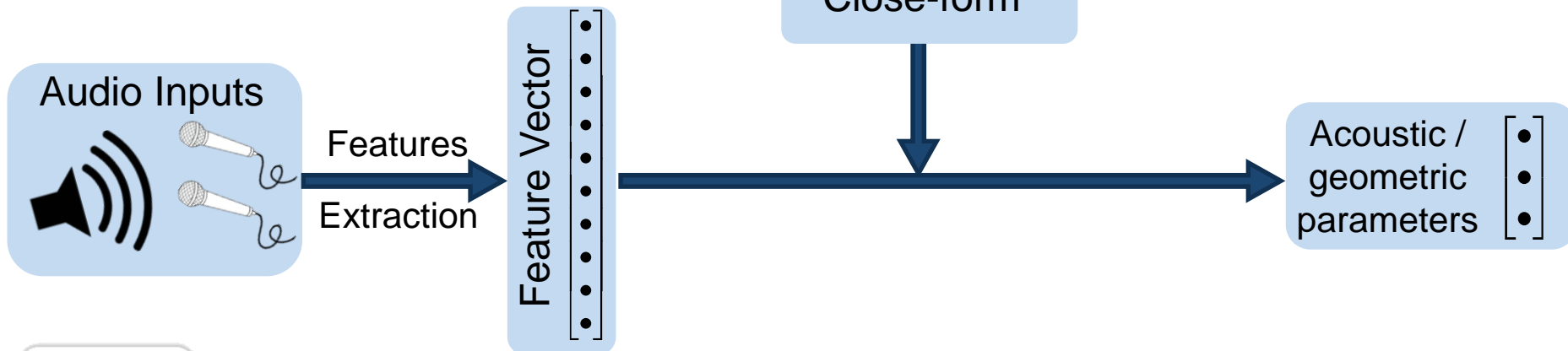
Present time, Inria Nancy

a) Physics-Driven / Traditional Approaches



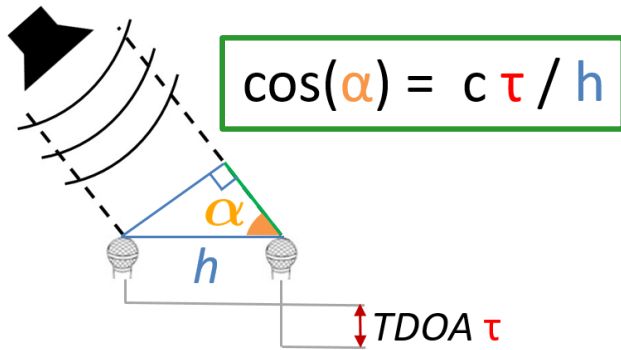
Sabine's law:

$$RT_{60}(b) \approx 0.16 \frac{V}{S \bar{\alpha}(b)}$$



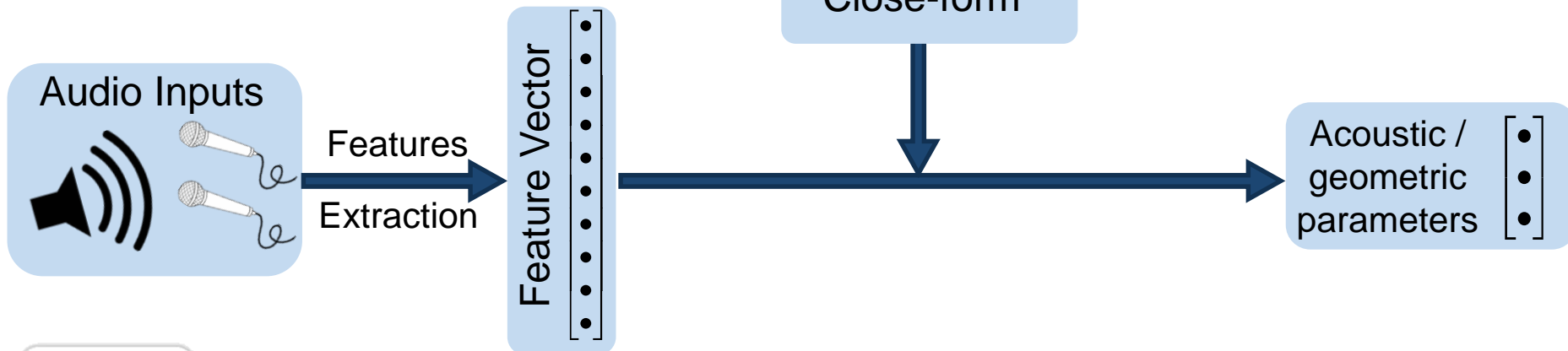
Present time, Inria Nancy

a) Physics-Driven / Traditional Approaches



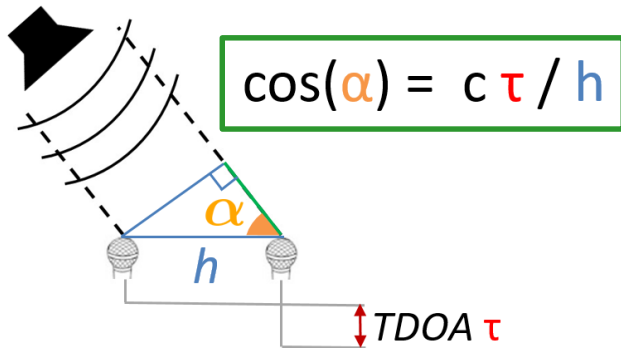
Sabine's law:

$$RT_{60}(b) \approx 0.16 \frac{V}{S \bar{\alpha}(b)}$$



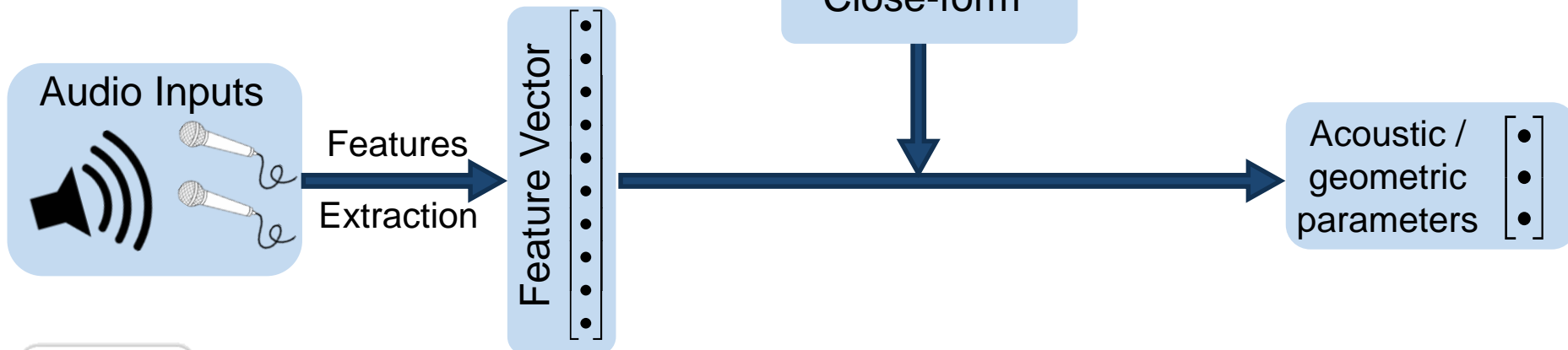
Present time, Inria Nancy

a) Physics-Driven / Traditional Approaches



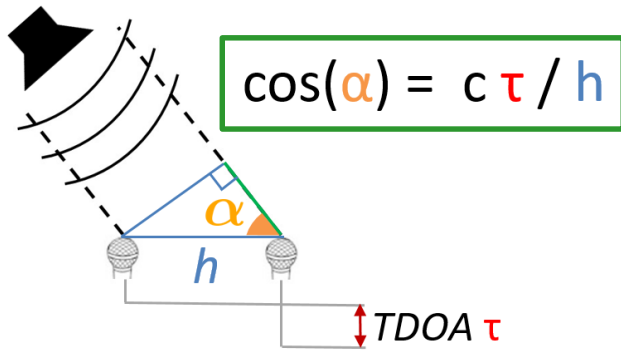
Sabine's law:

$$RT_{60}(b) \approx 0.16 \frac{V}{S \bar{\alpha}(b)}$$



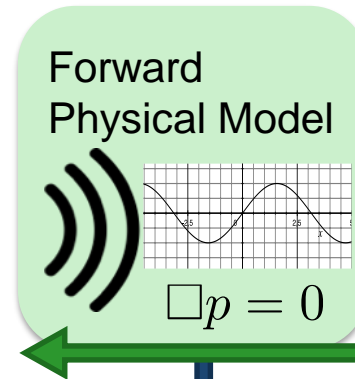
Present time, Inria Nancy

a) Physics-Driven / Traditional Approaches



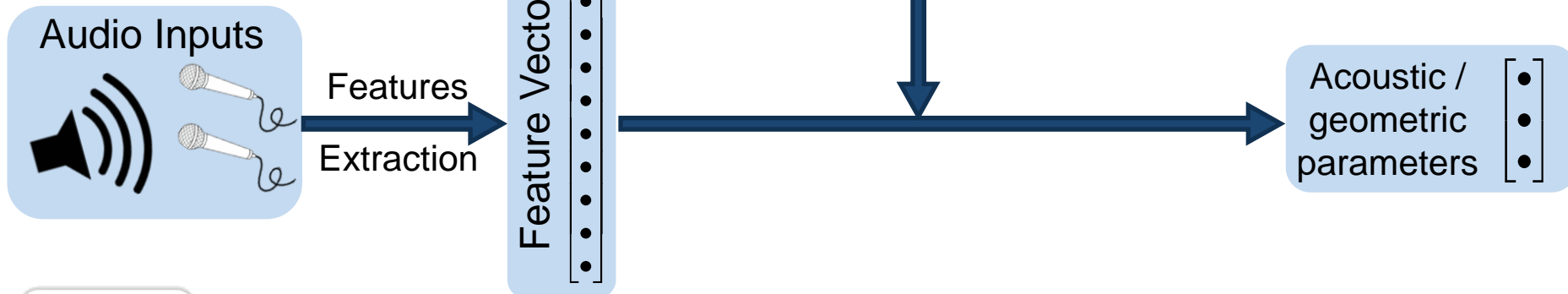
Sabine's law:

$$RT_{60}(b) \approx 0.16 \frac{V}{S \bar{\alpha}(b)}$$

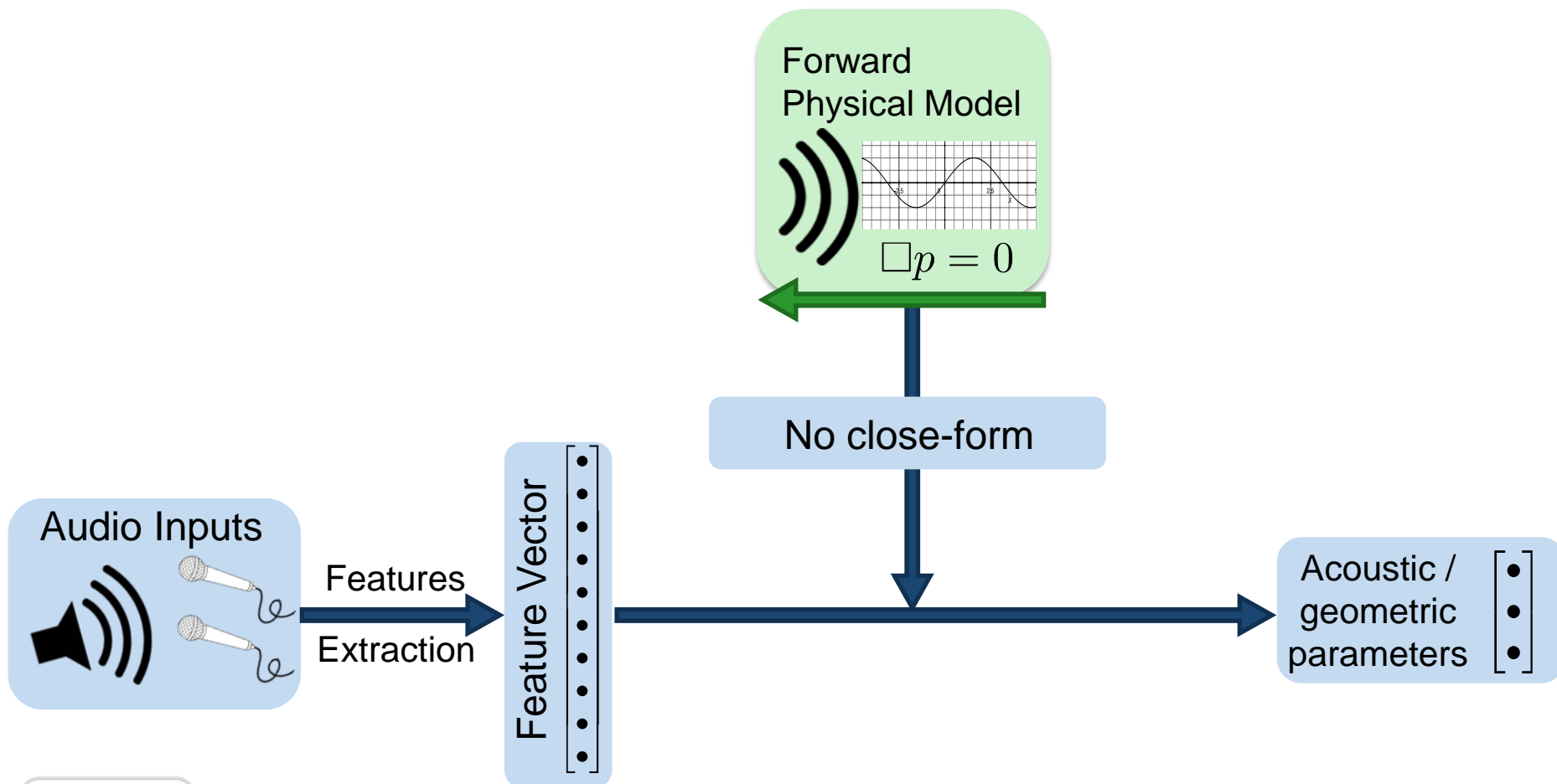


- ✓ No training data needed
- ✓ Computationally efficient
- ✗ Suffers in complex conditions
- ✗ Limited

Close-form



a) Physics-Driven / Traditional Approaches

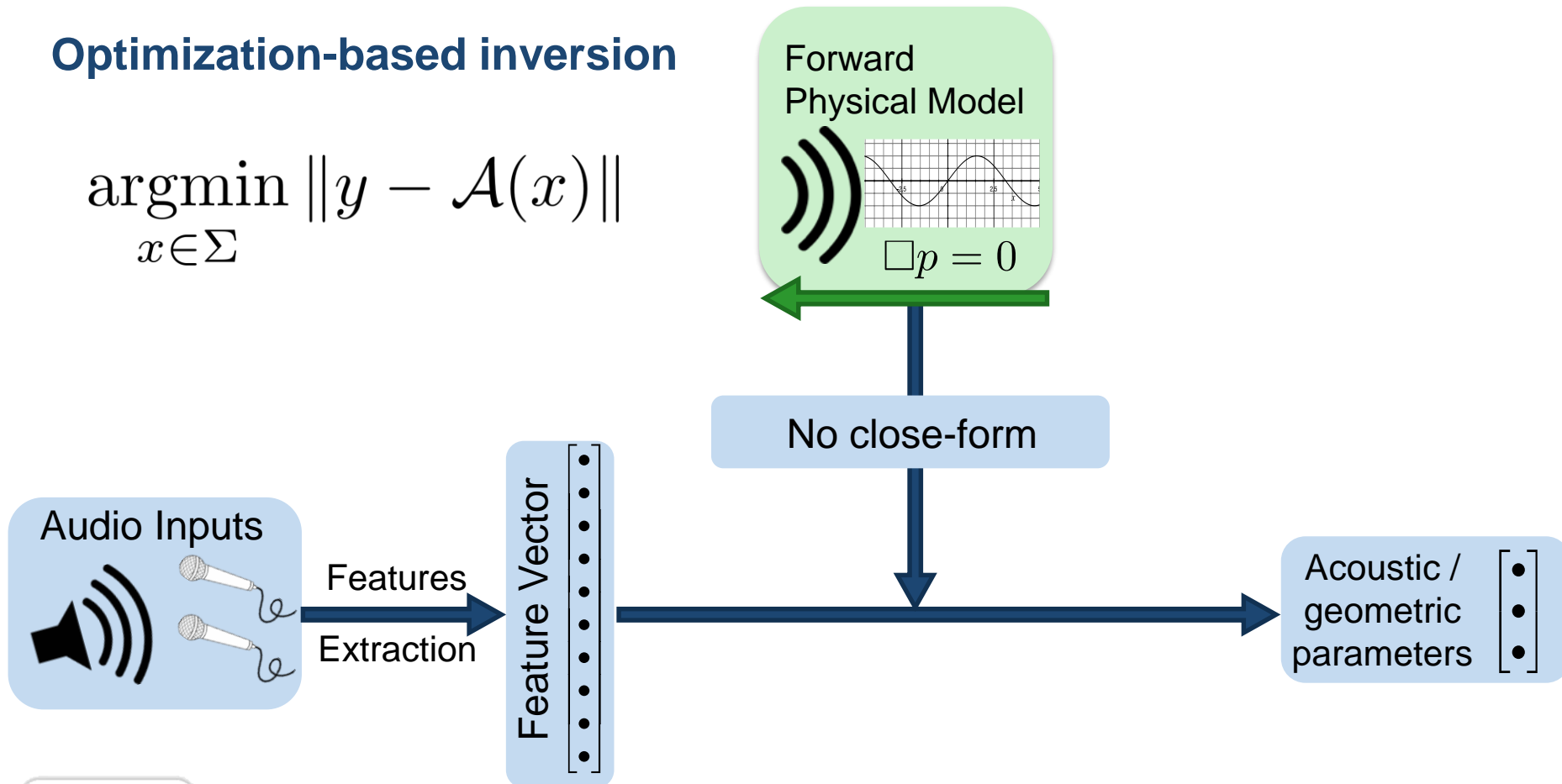


Present time, Inria Nancy

a) Physics-Driven / Traditional Approaches

Optimization-based inversion

$$\operatorname{argmin}_{x \in \Sigma} \|y - \mathcal{A}(x)\|$$

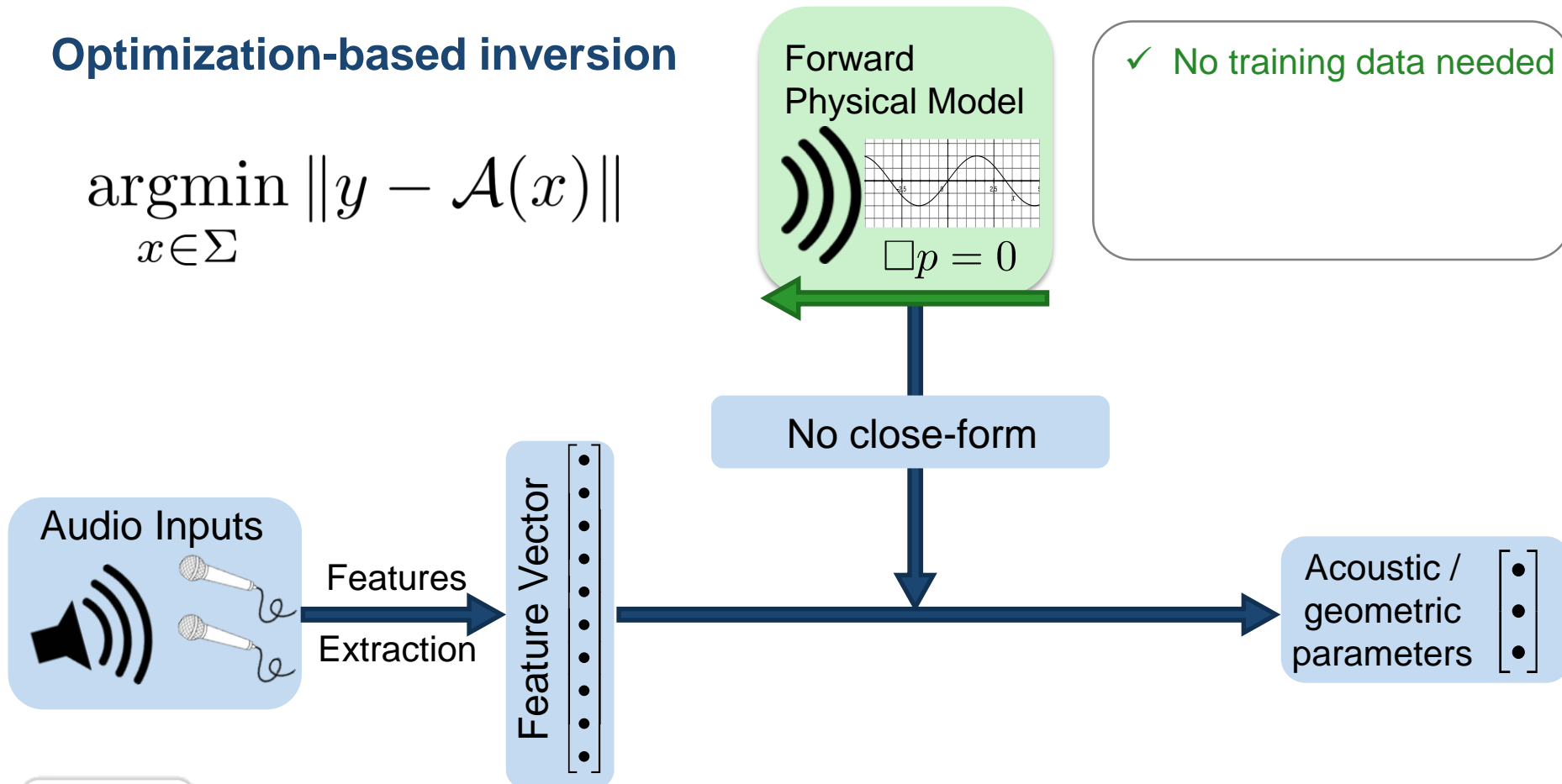


Present time, Inria Nancy

a) Physics-Driven / Traditional Approaches

Optimization-based inversion

$$\operatorname{argmin}_{x \in \Sigma} \|y - \mathcal{A}(x)\|$$

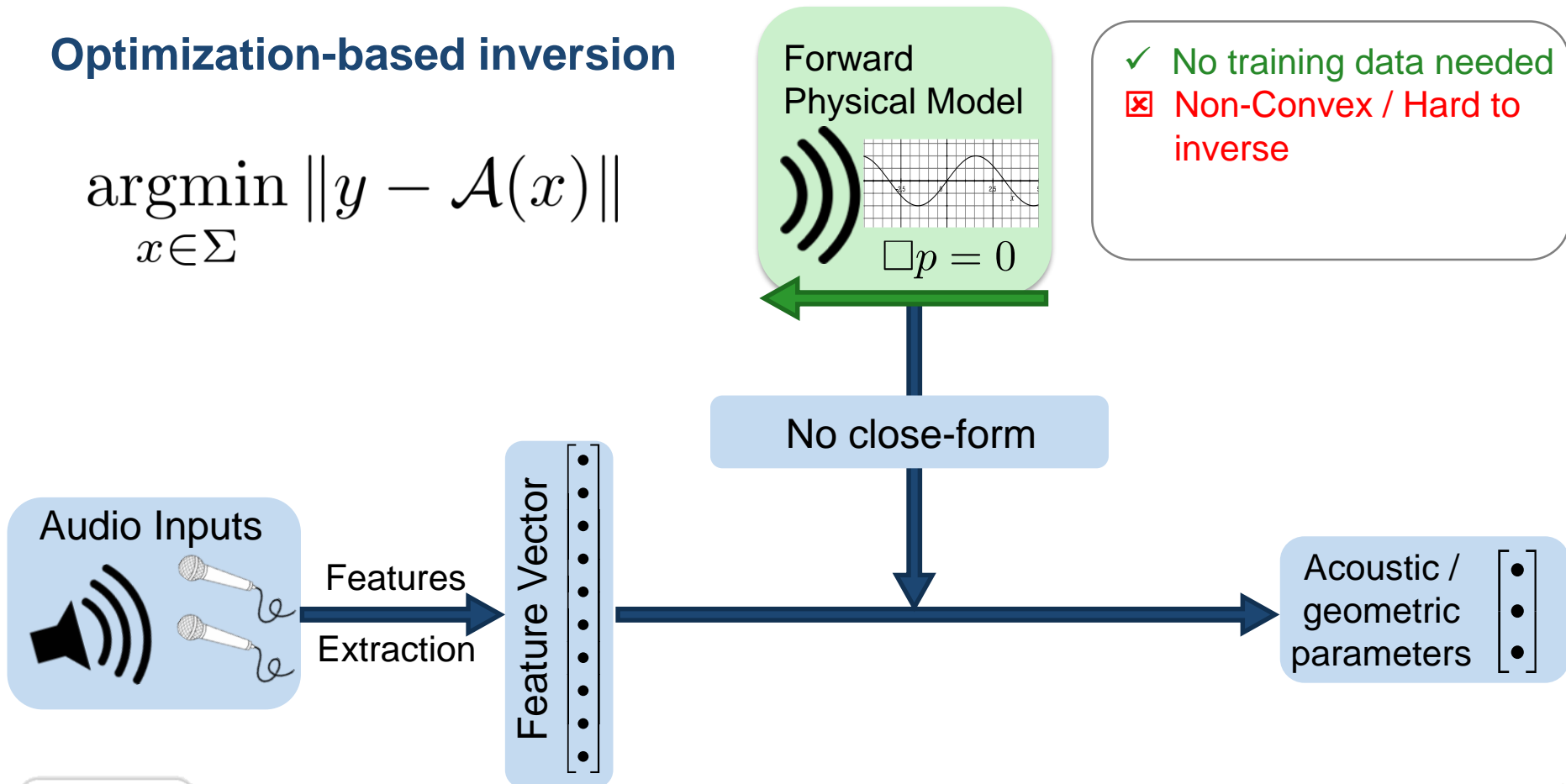


Present time, Inria Nancy

a) Physics-Driven / Traditional Approaches

Optimization-based inversion

$$\operatorname{argmin}_{x \in \Sigma} \|y - \mathcal{A}(x)\|$$

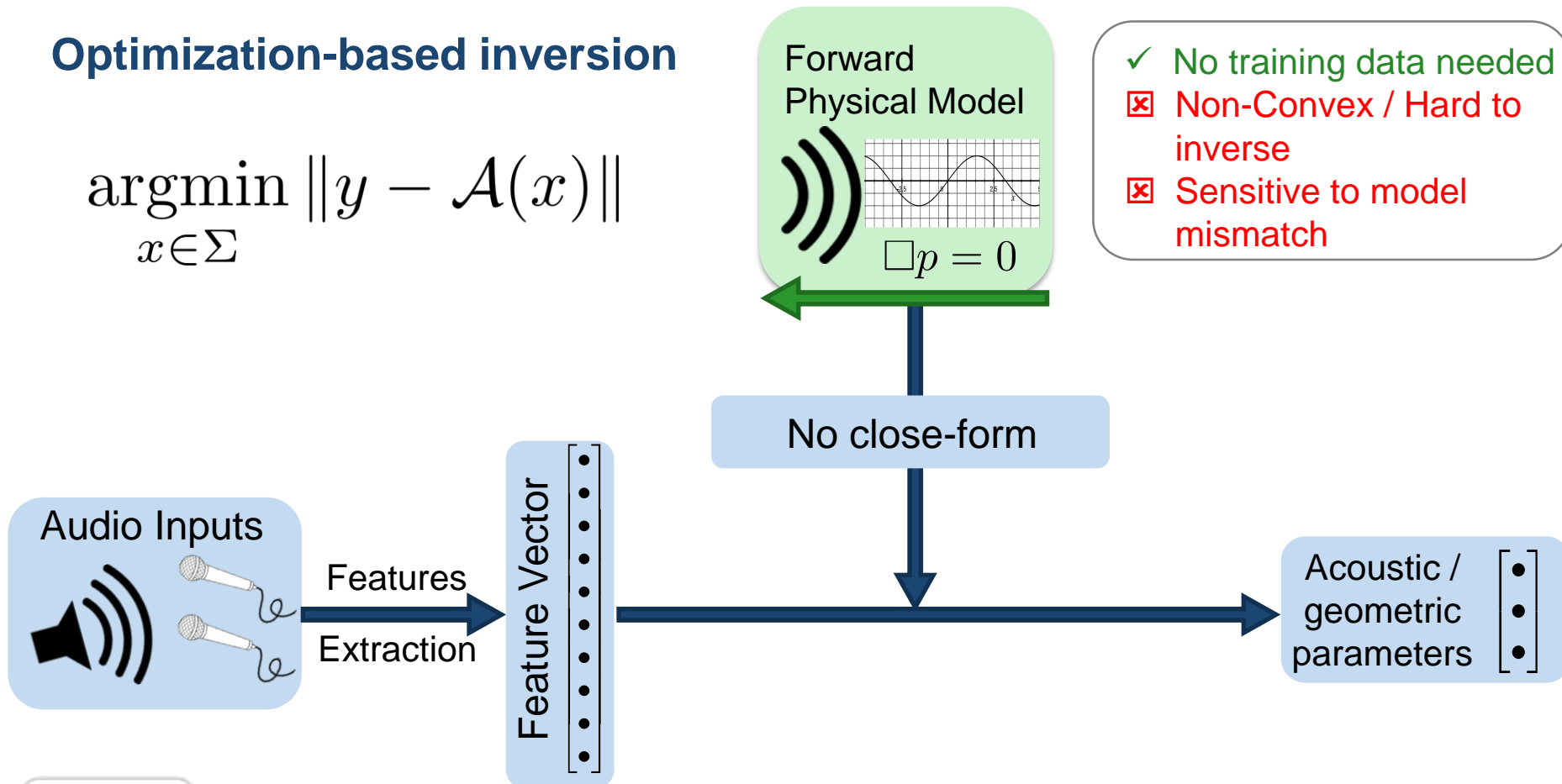


Present time, Inria Nancy

a) Physics-Driven / Traditional Approaches

Optimization-based inversion

$$\operatorname{argmin}_{x \in \Sigma} \|y - \mathcal{A}(x)\|$$

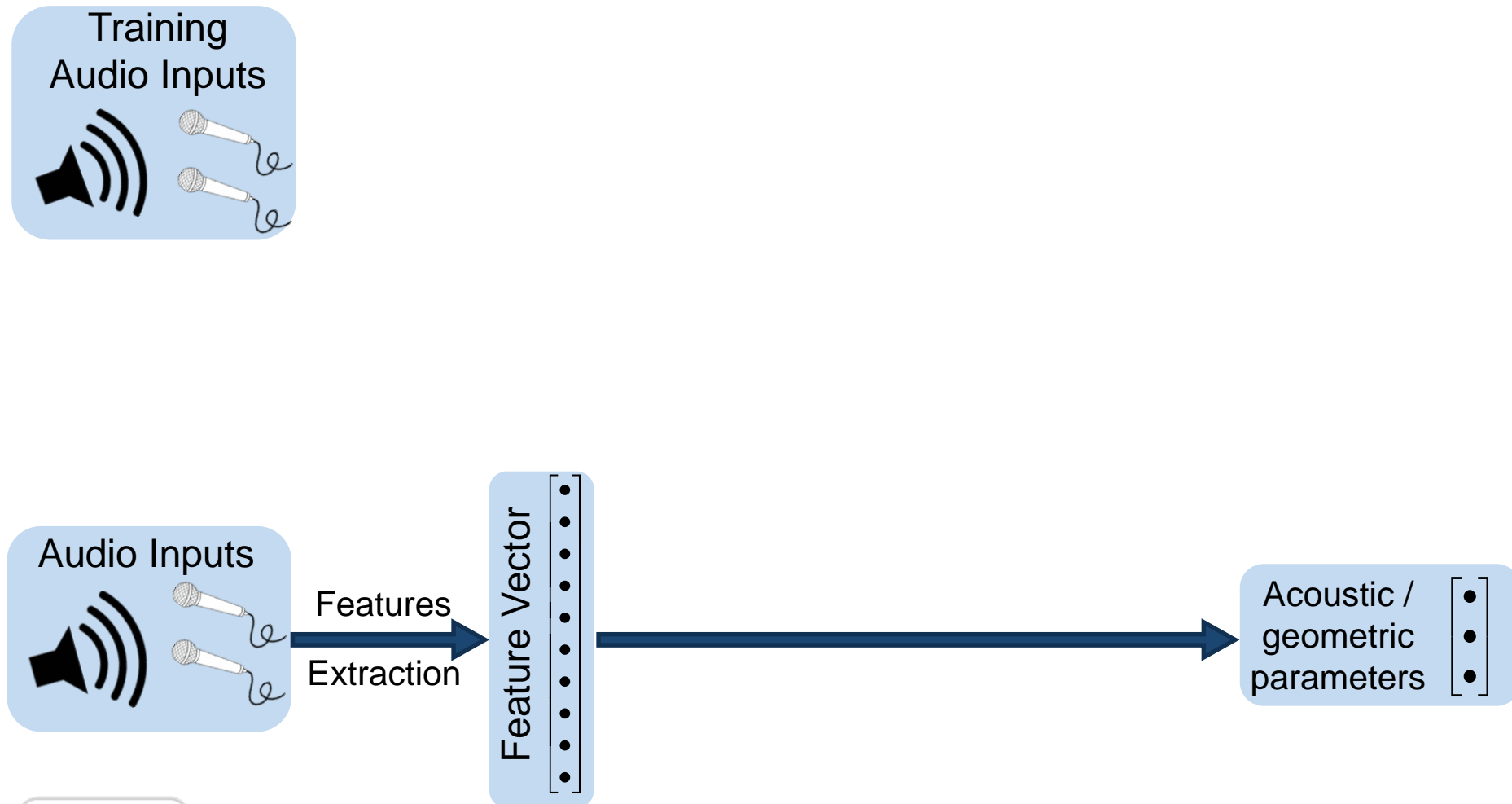


b) Real-Data-Driven Approaches



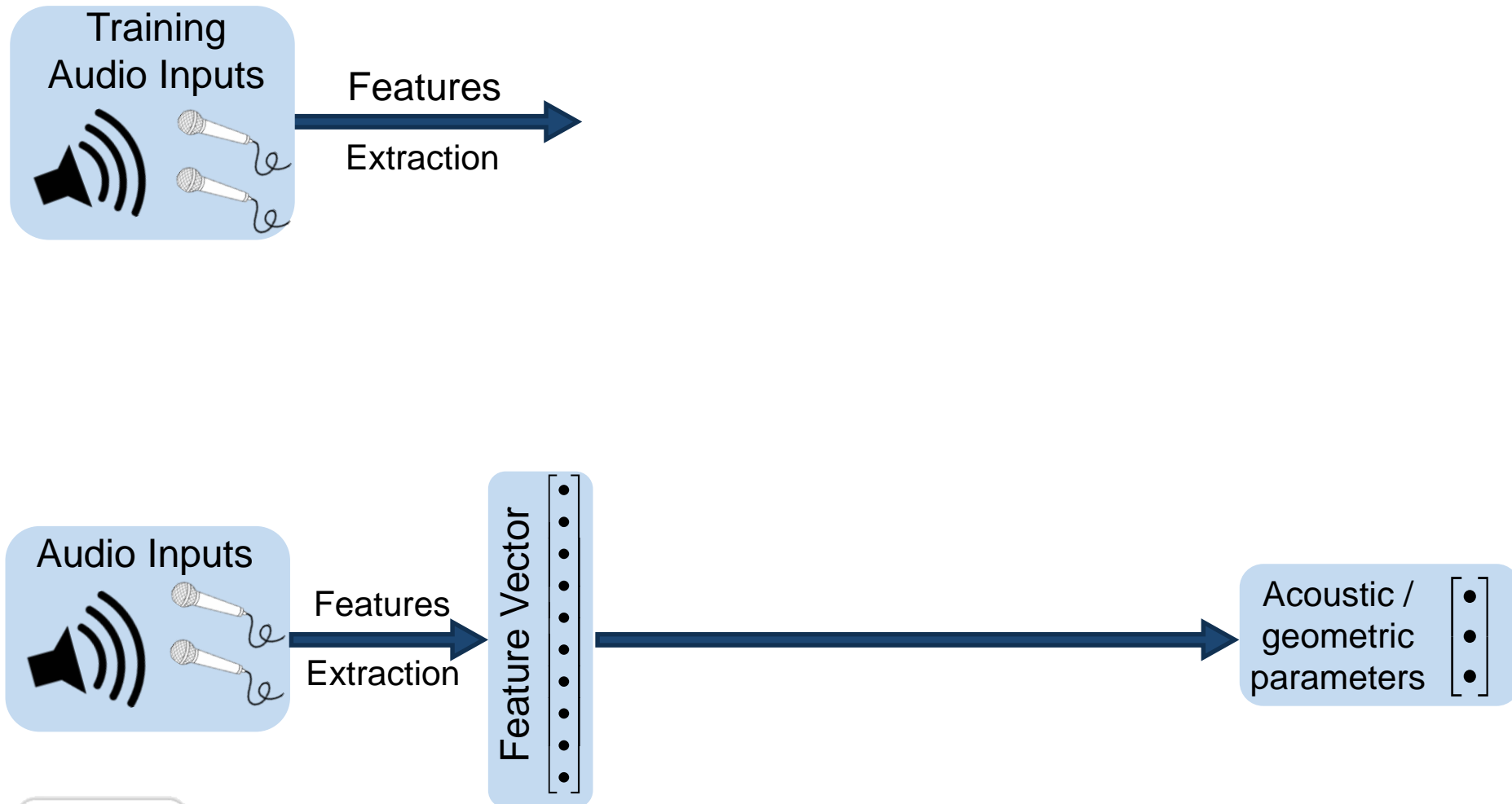
Present time, Inria Nancy

b) Real-Data-Driven Approaches



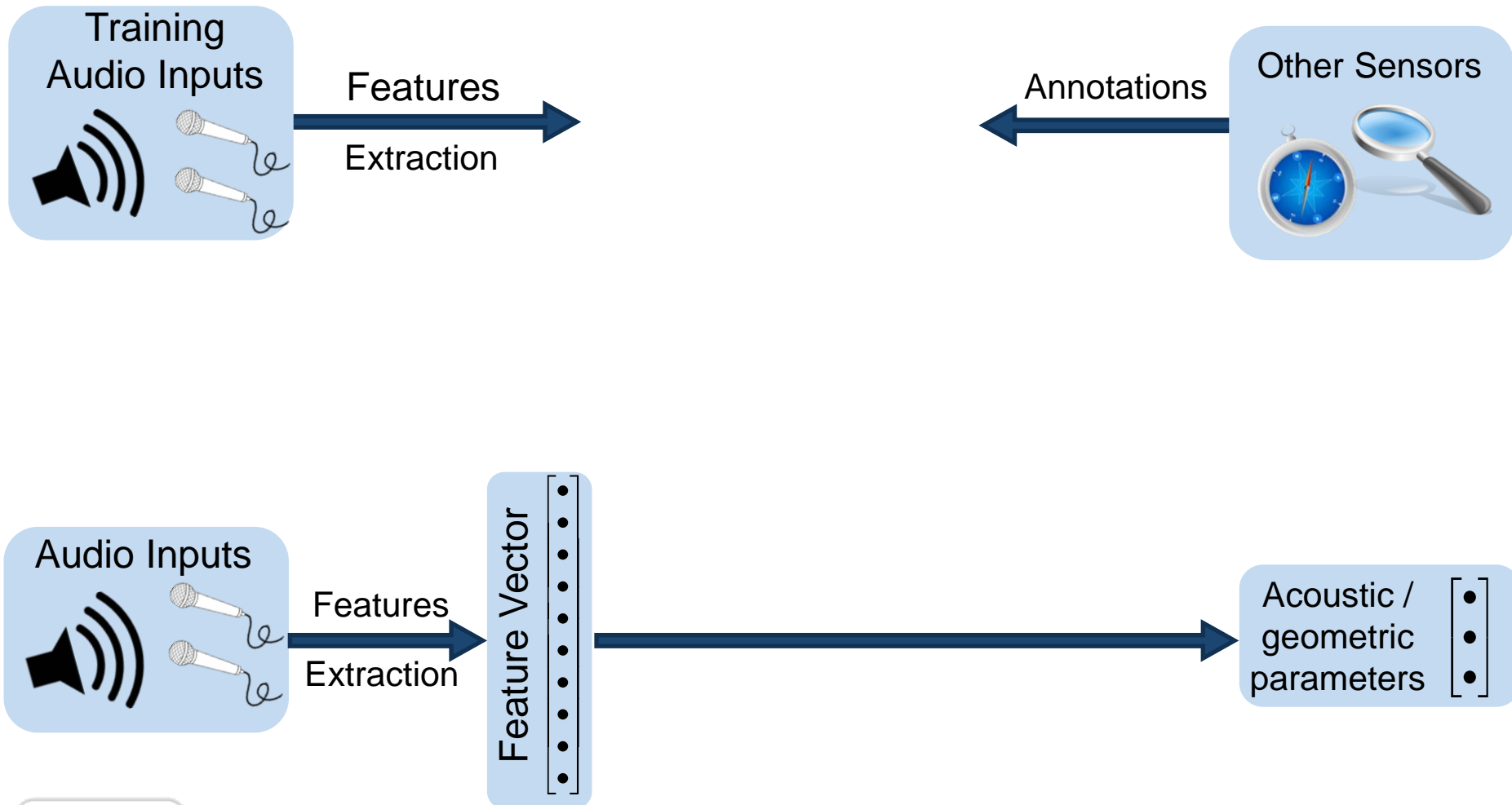
Present time, Inria Nancy

b) Real-Data-Driven Approaches



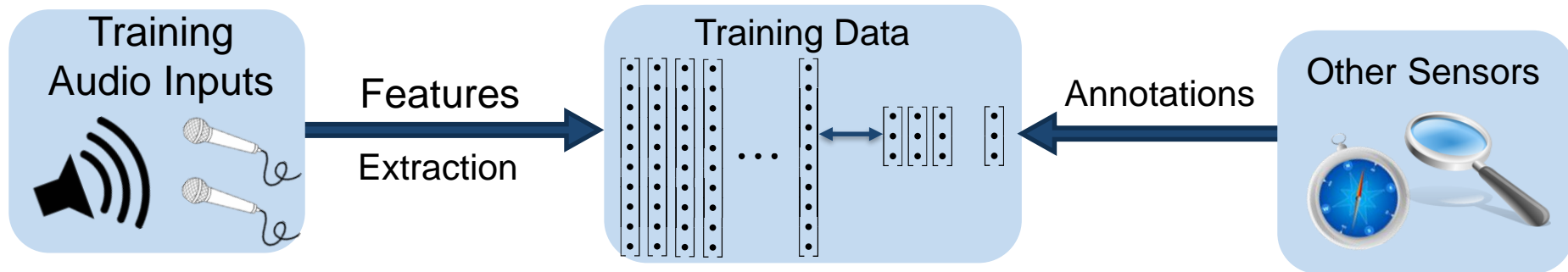
Present time, Inria Nancy

b) Real-Data-Driven Approaches



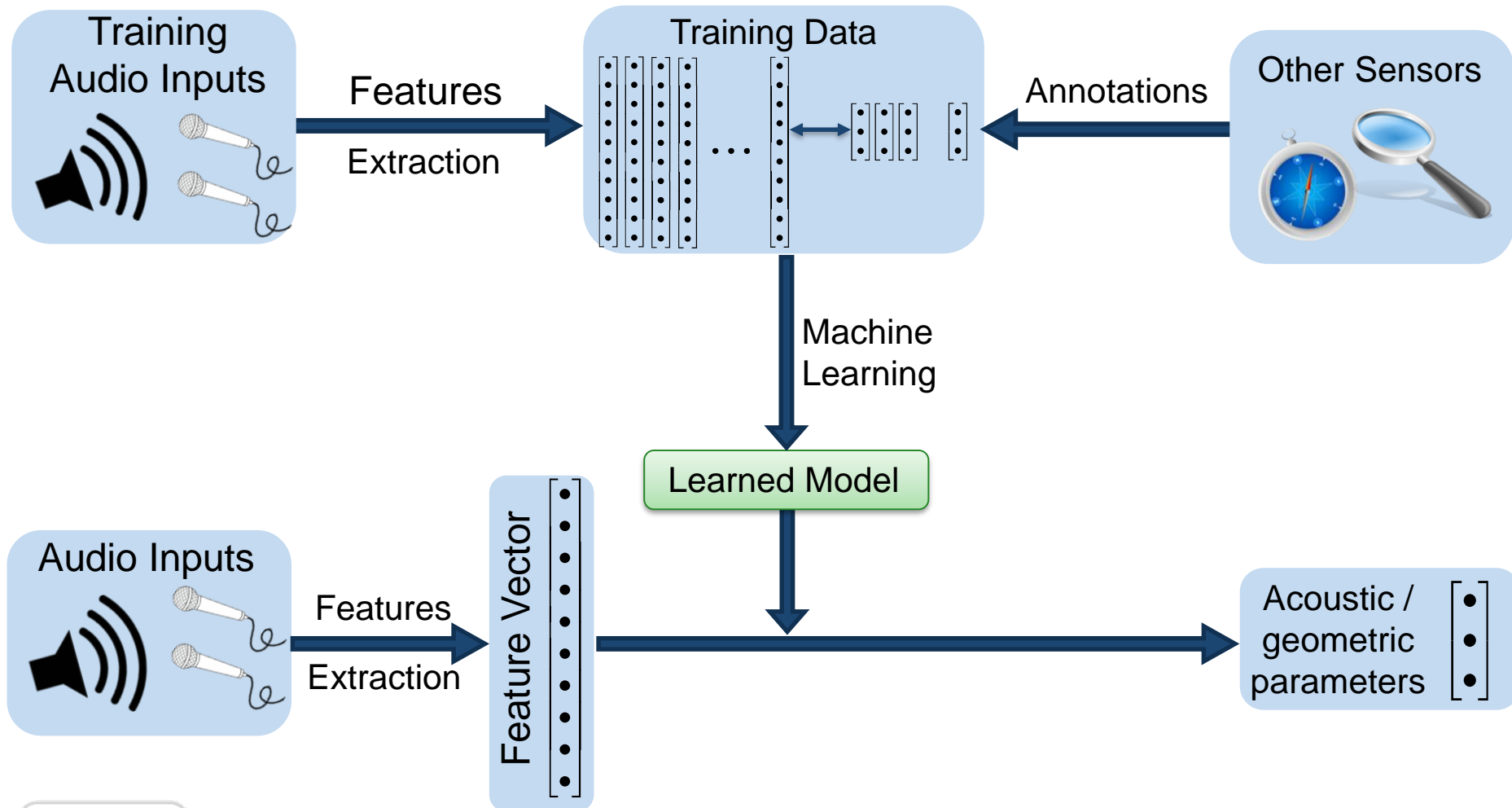
Present time, Inria Nancy

b) Real-Data-Driven Approaches



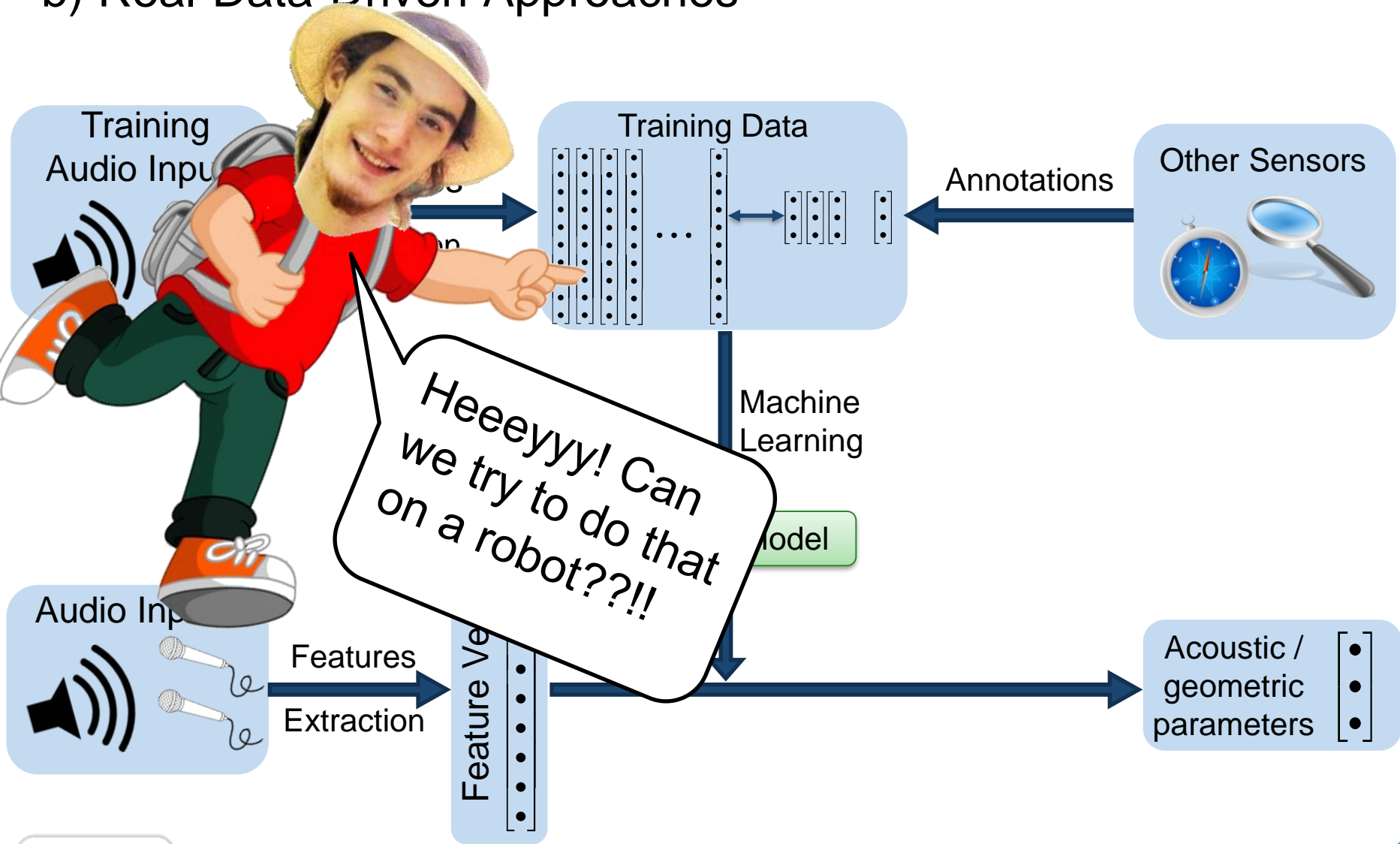
Present time, Inria Nancy

b) Real-Data-Driven Approaches



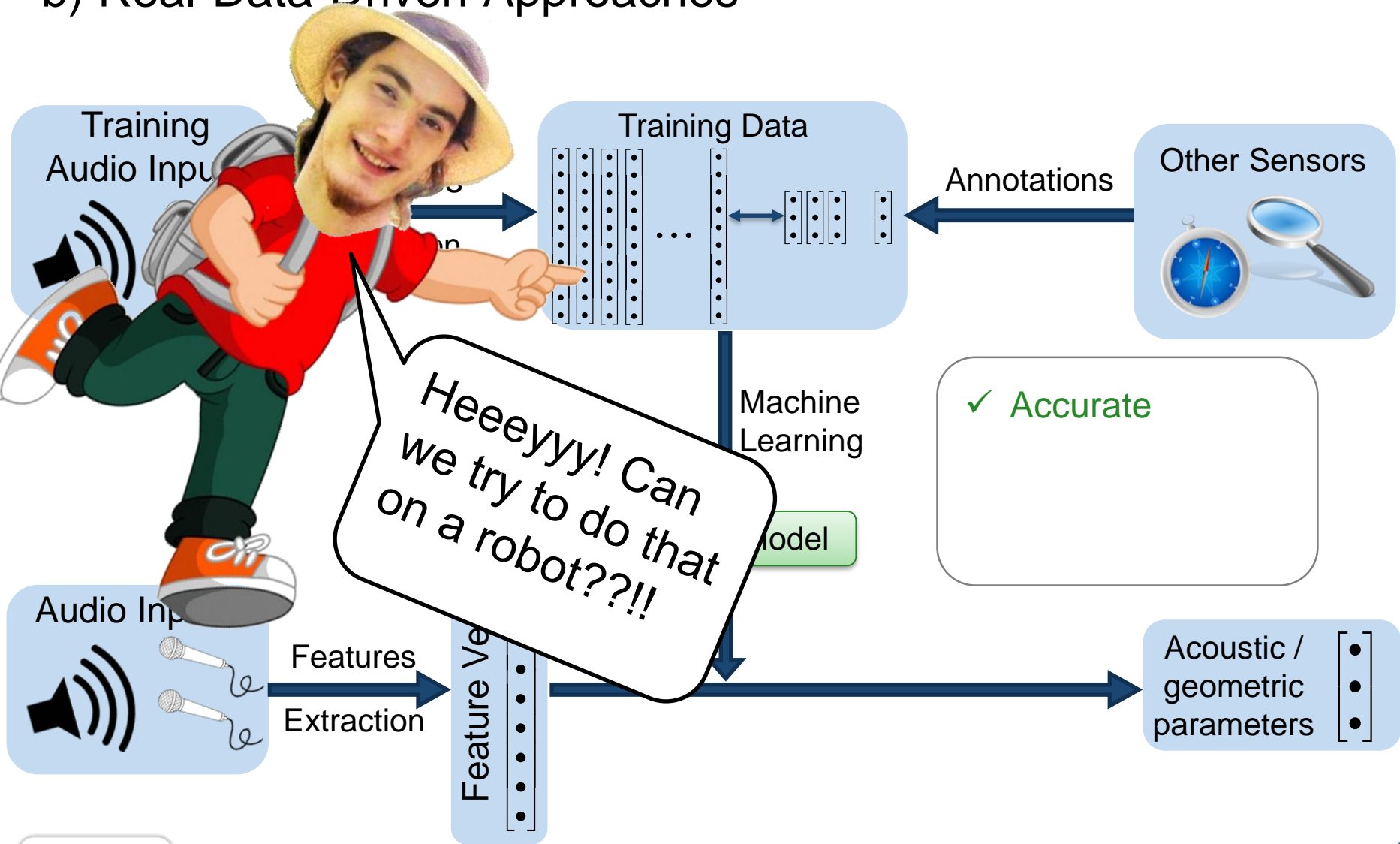
Present time, Inria Nancy

b) Real-Data-Driven Approaches



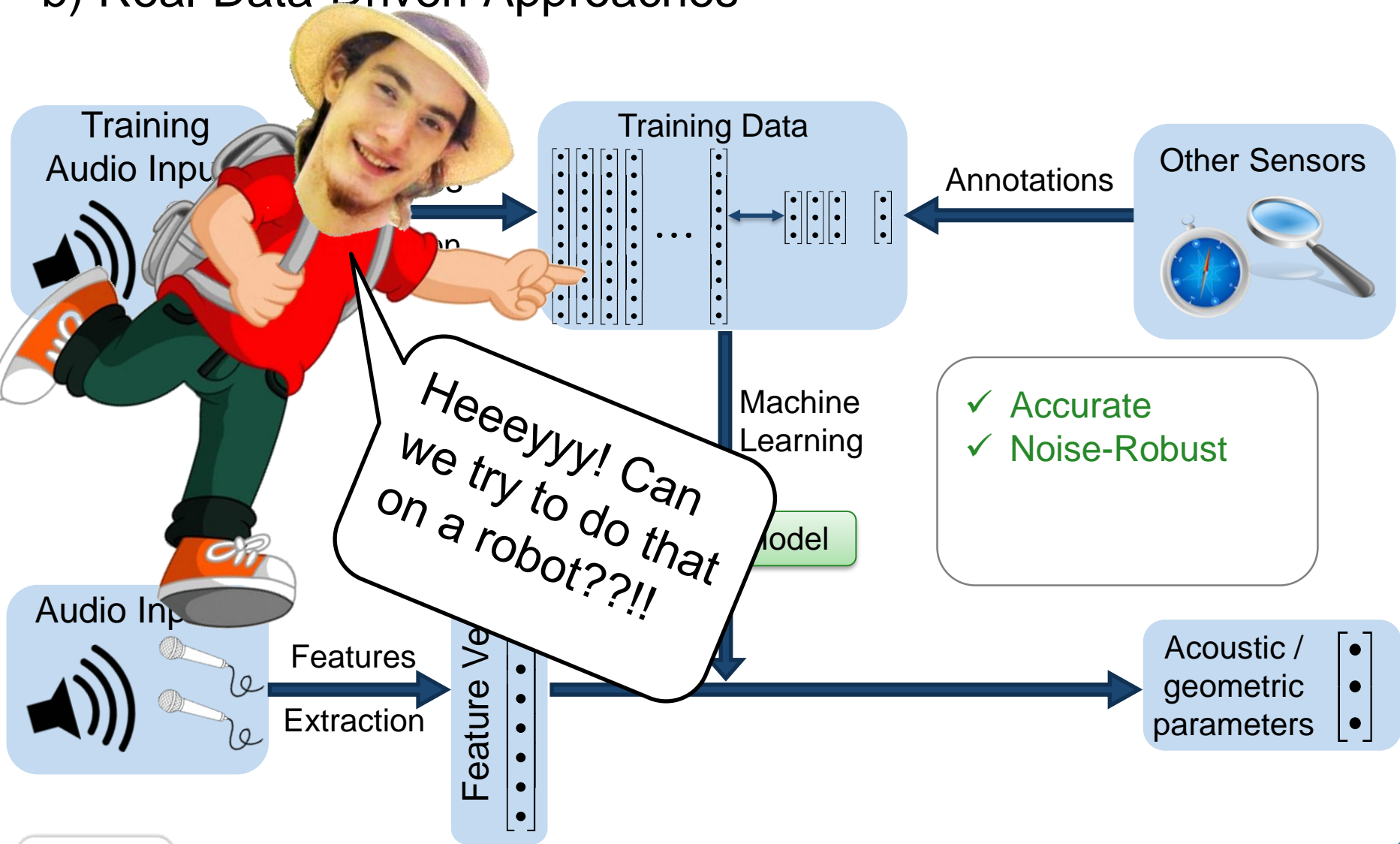
Present time, Inria Nancy

b) Real-Data-Driven Approaches



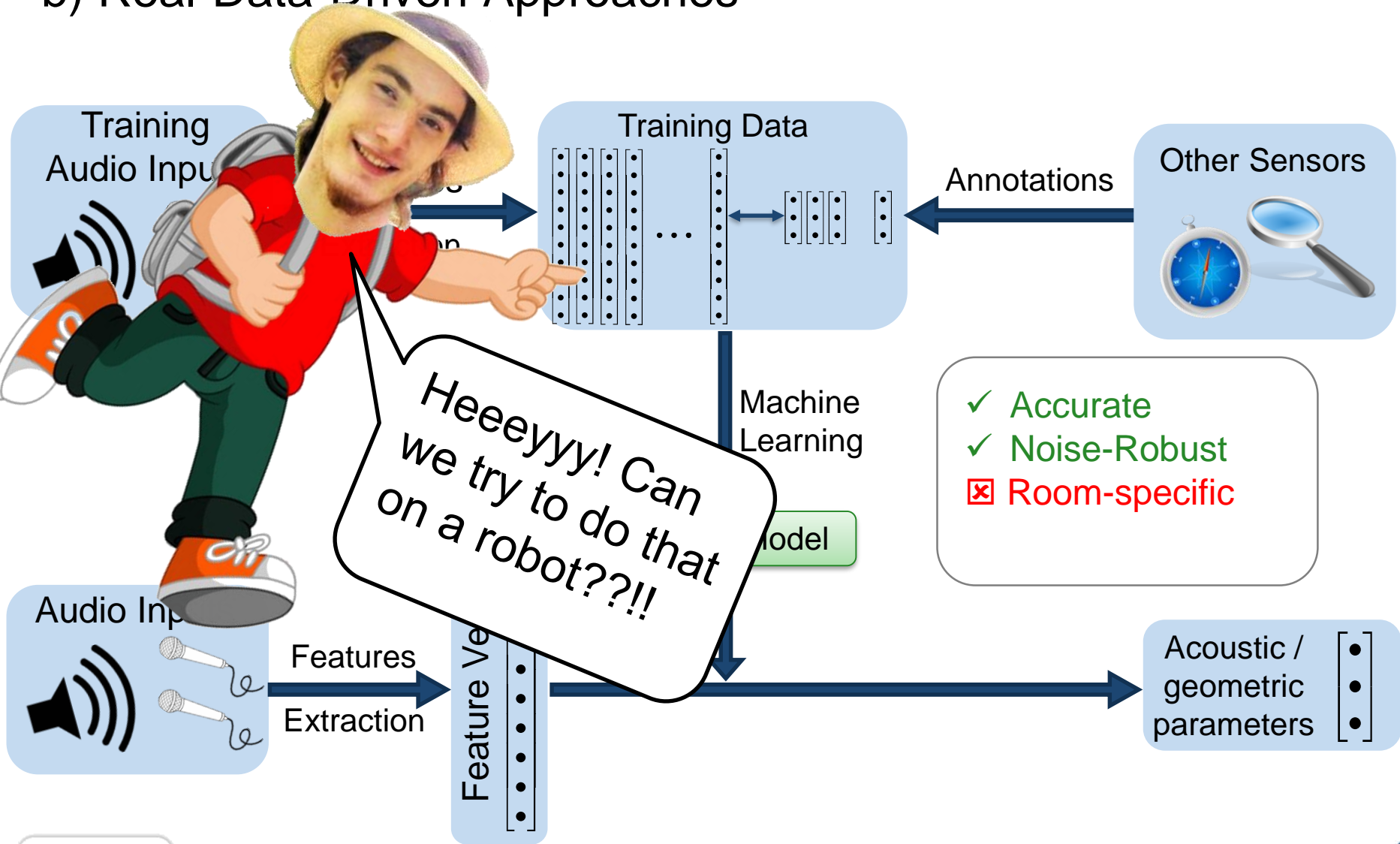
Present time, Inria Nancy

b) Real-Data-Driven Approaches



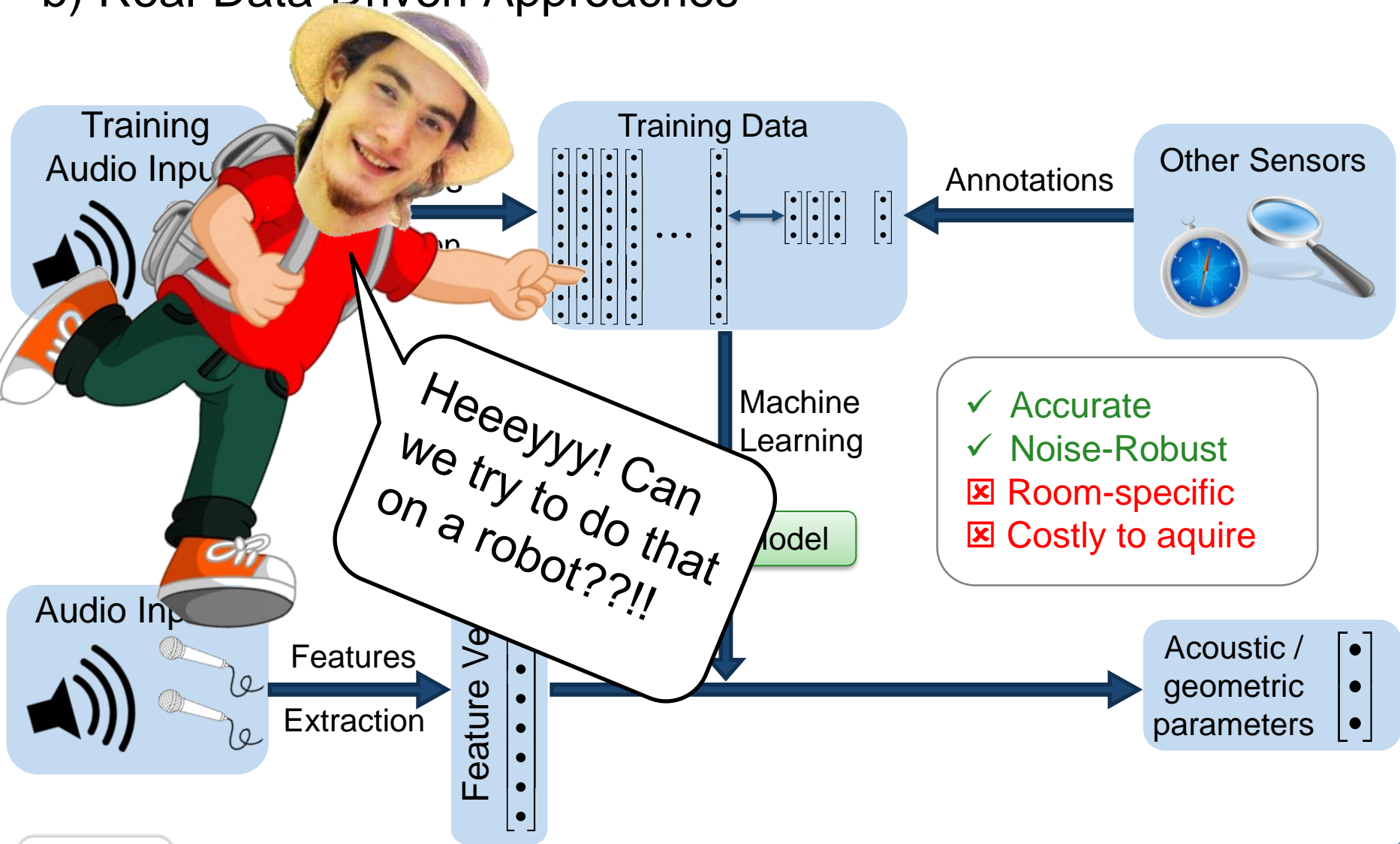
Present time, Inria Nancy

b) Real-Data-Driven Approaches



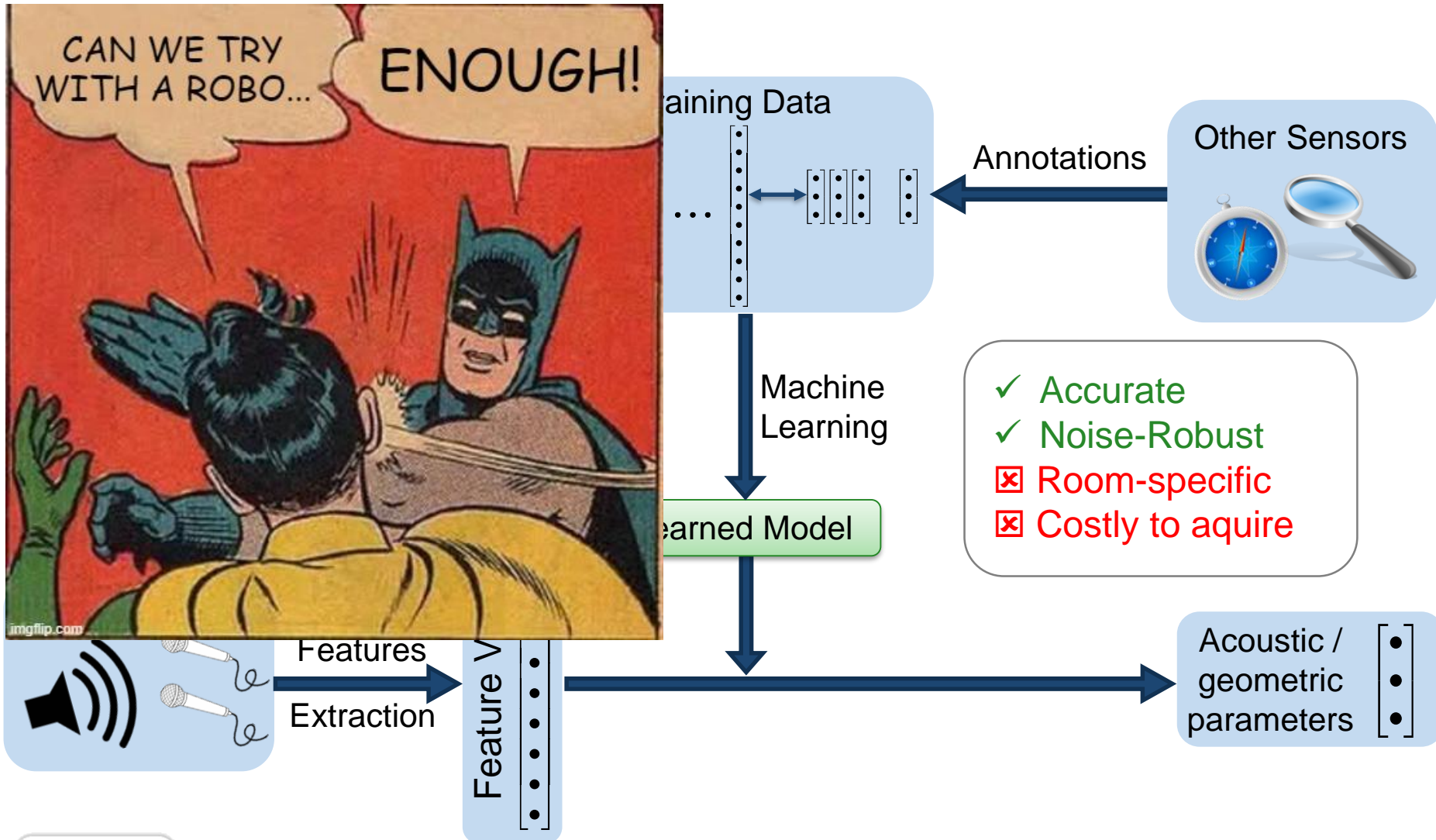
Present time, Inria Nancy

b) Real-Data-Driven Approaches



Present time, Inria Nancy

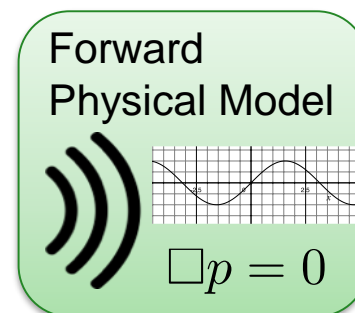
b) Real-Data-Driven Approaches



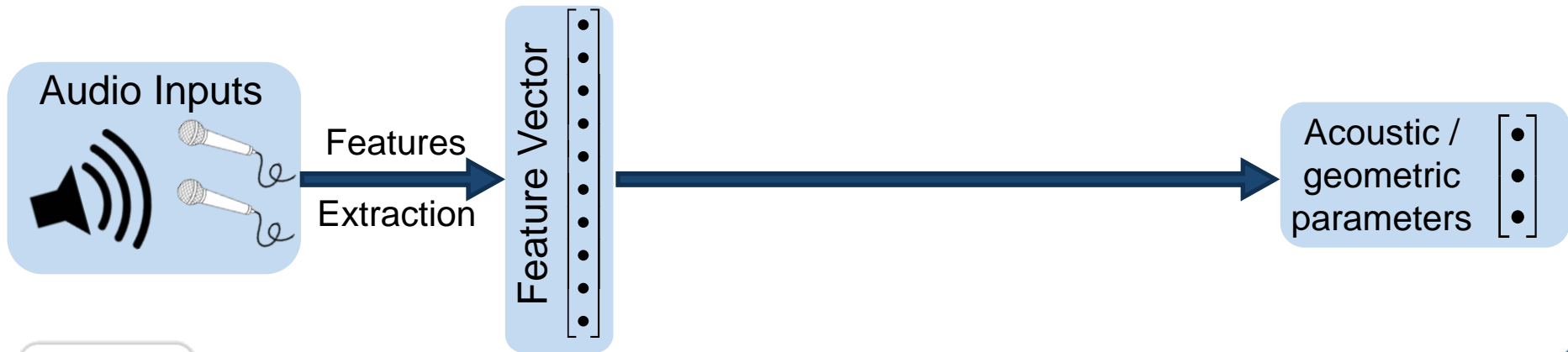
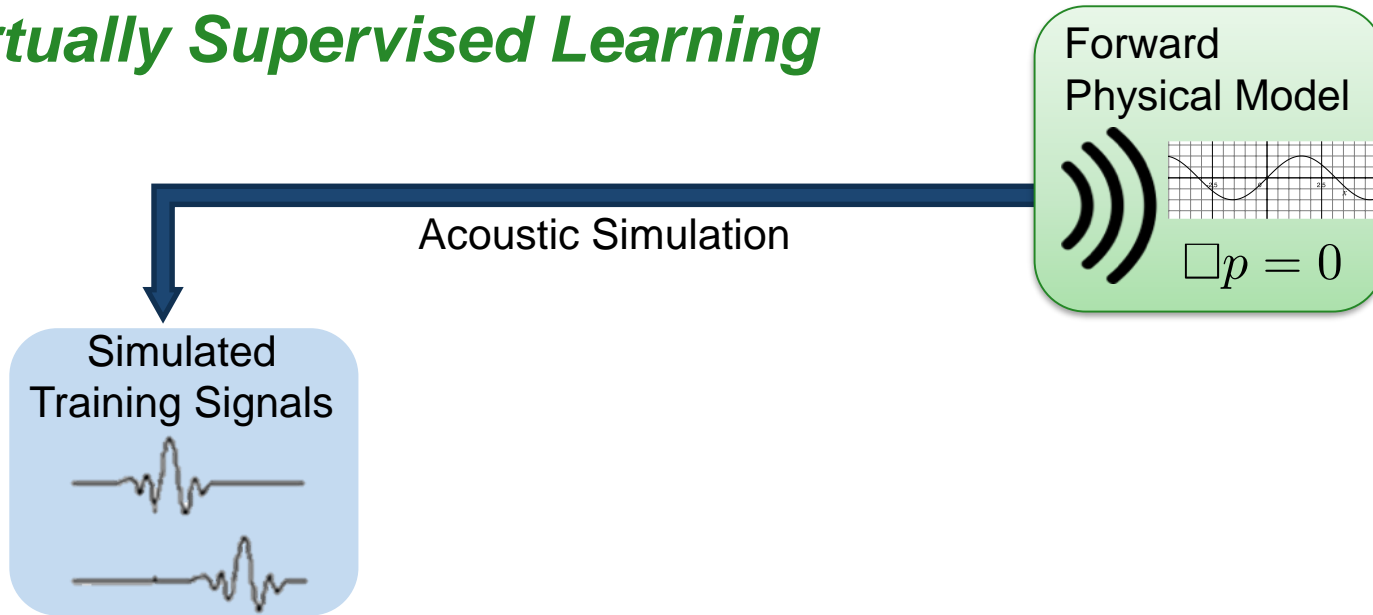
c) *Virtually Supervised Learning*



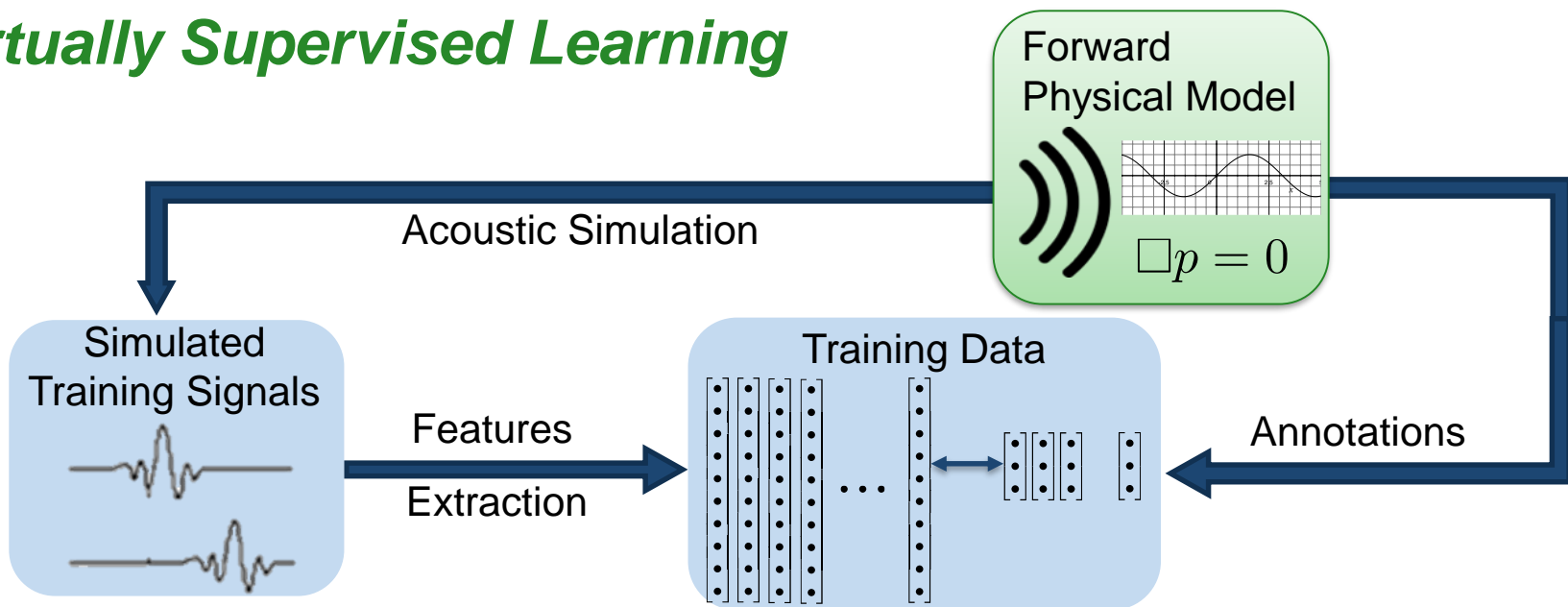
c) *Virtually Supervised Learning*



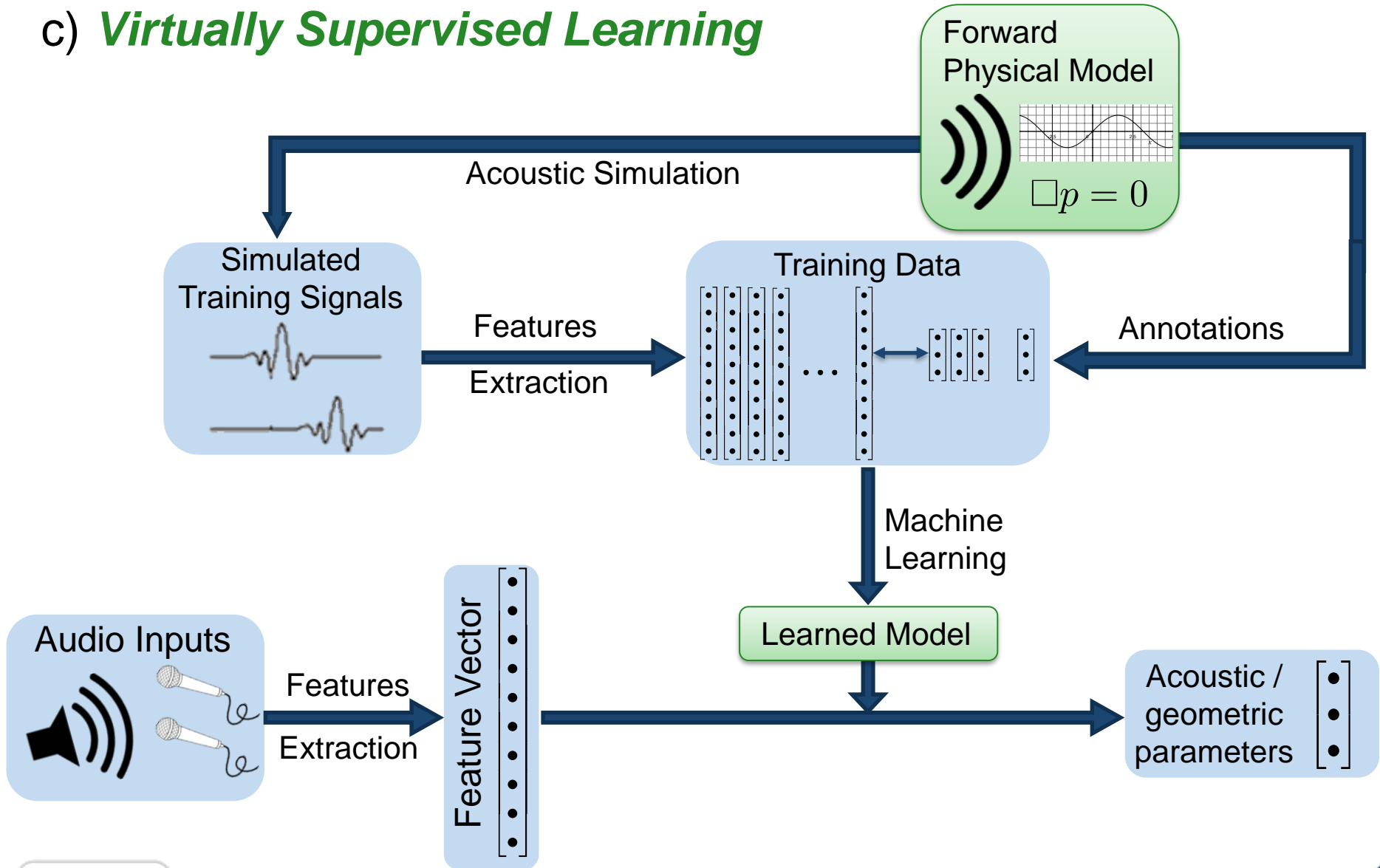
c) *Virtually Supervised Learning*



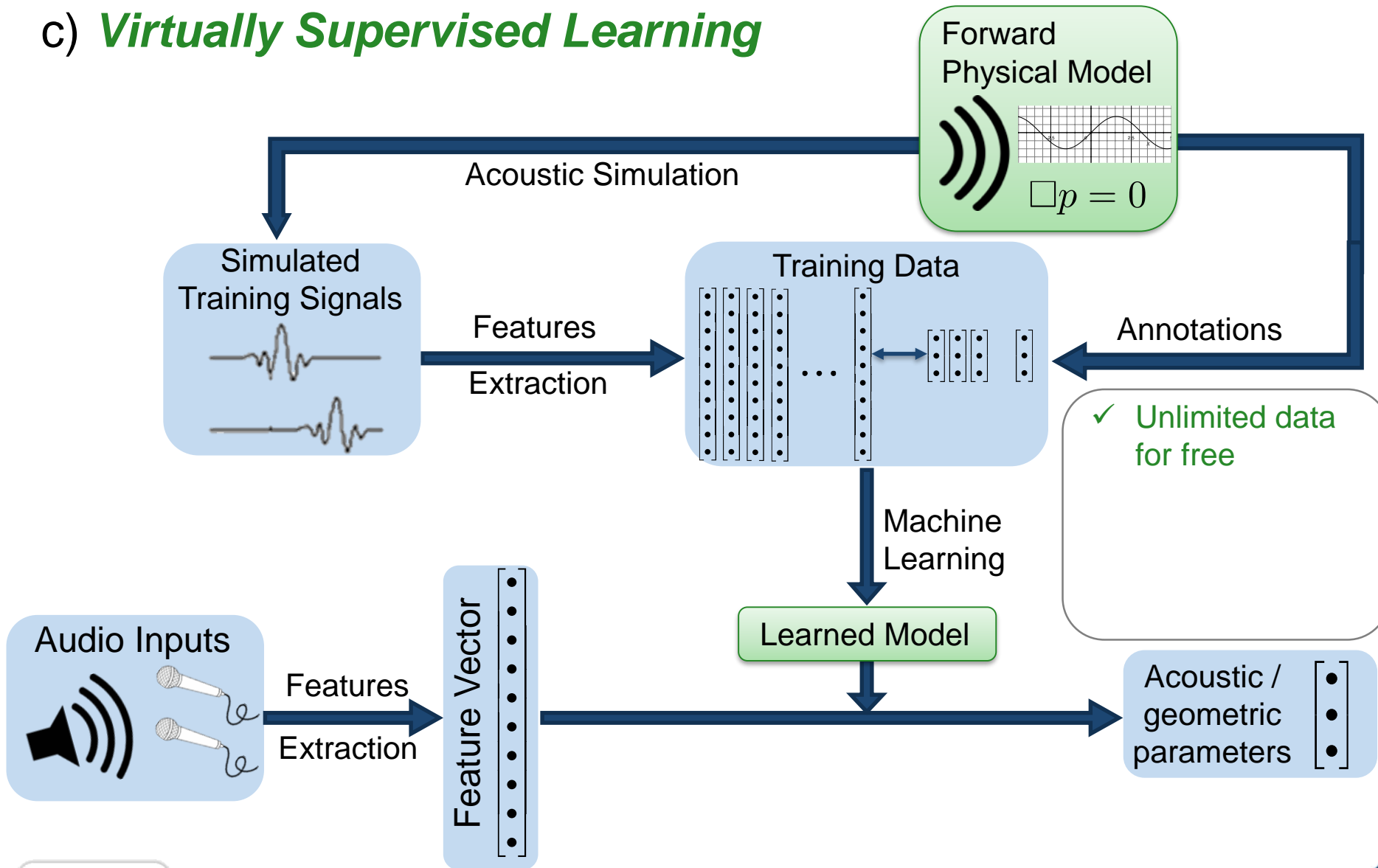
c) *Virtually Supervised Learning*



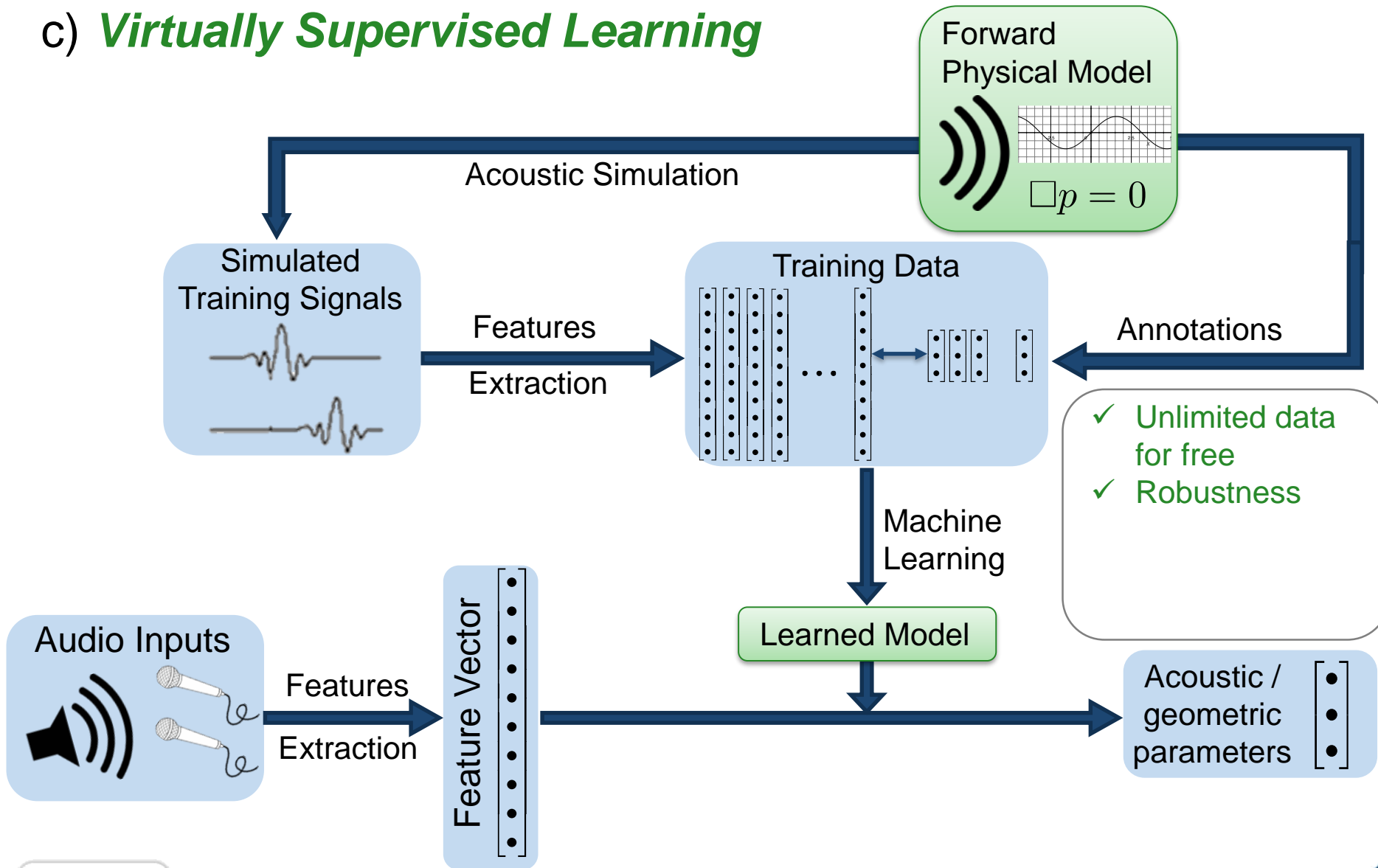
c) *Virtually Supervised Learning*



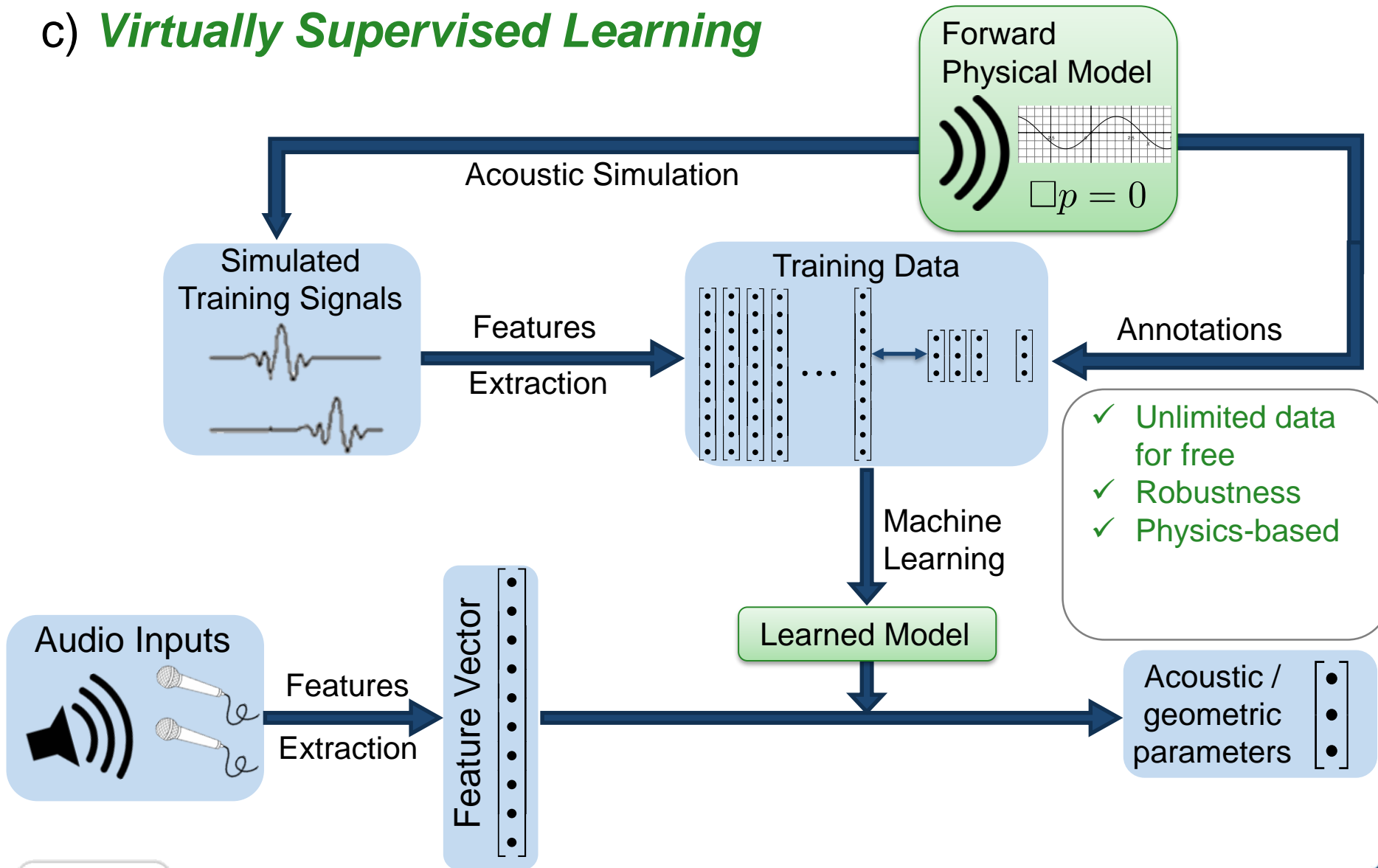
c) *Virtually Supervised Learning*



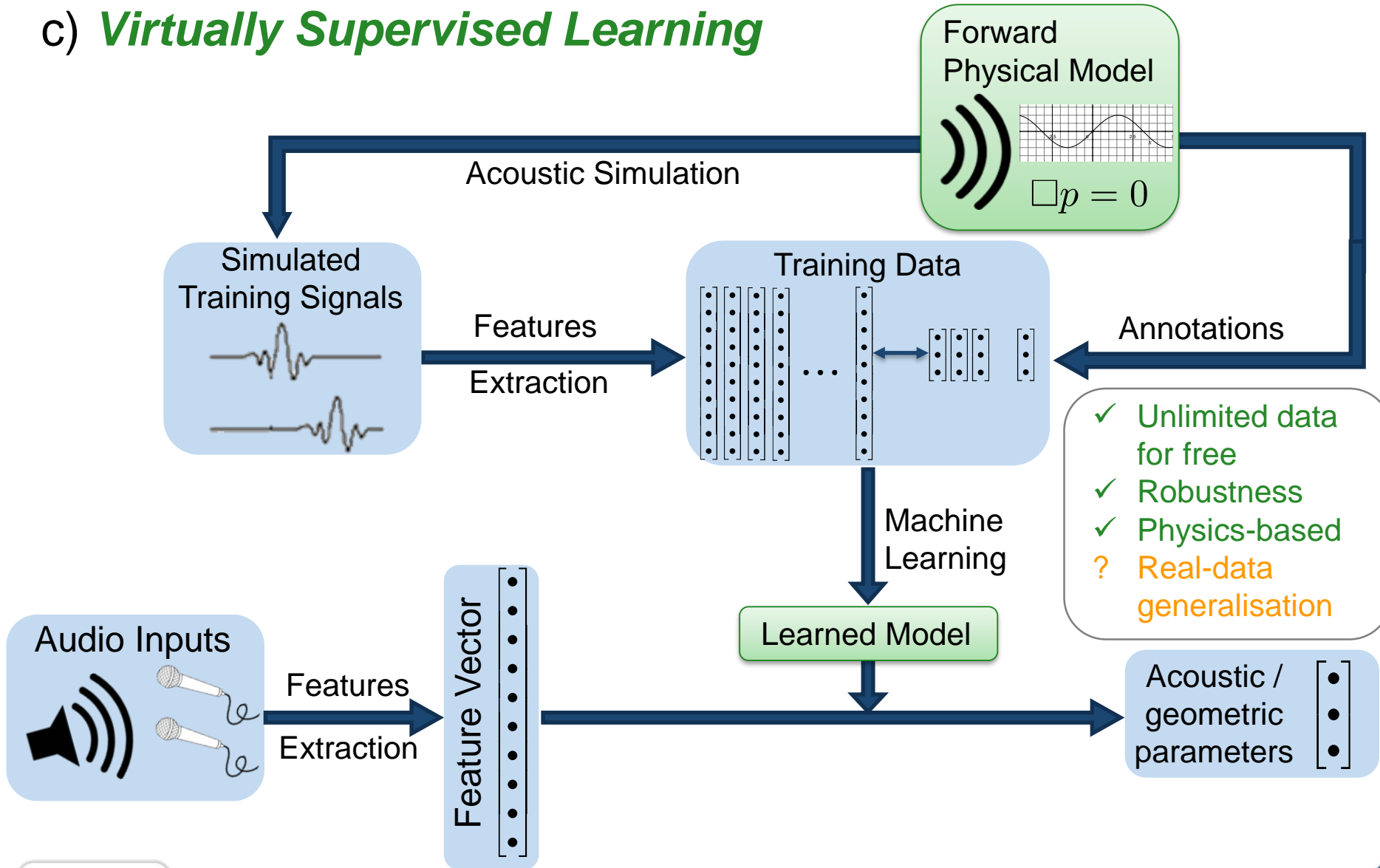
c) *Virtually Supervised Learning*



c) *Virtually Supervised Learning*



c) *Virtually Supervised Learning*



HOW TO (VIRTUALLY) TRAIN YOUR SOUND SOURCE LOCALIZER

Prerak Srivastava¹, Antoine Deleforge¹, Archontis Politis², Emmanuel Vincent¹



HOW TO (VIRTUALLY) TRAIN YOUR SOUND SOURCE LOCALIZER

Prerak Srivastava¹, Antoine Deleforge¹, Archontis Politis², Emmanuel Vincent¹

JASA REVIEW-ARTICLE

July 2022



A survey of sound source localization with deep learning methods

Pierre-Amaury Grumiaux,^{1,a)} Srđan Kitić,² Laurent Girin,³ and Alexandre Guérin²

¹Nantes Université, École Centrale Nantes, CNRS, LS2N, 2 chemin de la Houssinière, F-44332 Nantes, France

²Orange Labs, 4 Rue du Clos Courtel, 35510 Cesson-Sévigné, France

³Univ. Grenoble Alpes, Grenoble-INP, GIPSA-lab, 11 Rue des Mathématiques, 38400 Saint-Martin-d'Hères, France



HOW

Prerak

J.A.

A survey
method

Pierre-Ama

¹Nantes Unive

²Orange Labs,

³Univ. Grenob

| Author | Year | Architecture | Type | Learn. | Input features | Output | Sources | | | Data | | | | | | | |
|--------------------------------------|------|--------------------|------|--------|---------------------------------------------|----------------|----------|------|------|-------|----|----|----|------|----|----|----|
| | | | | | | | NoS | Kno. | Mov. | Train | | | | Test | | | |
| | | | | | | | | | | SA | RA | SR | RR | SA | RA | SR | RR |
| Adavanne <i>et al.</i> (2021) | 2021 | CRNN + SA | R | S | FOA Mel spectrograms + intensity + GCC-PHAT | x, y, z | 2 | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Bai <i>et al.</i> (2021) | 2021 | Res. CRNN | R | S | Log-Mel spectrograms + intensity | x, y, z | 1 | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Bianco <i>et al.</i> (2021) | 2021 | VAE | C | SS | RTF | θ | 1 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Bohlender <i>et al.</i> (2021) | 2021 | CNN/CRNN | C | S | Phase map | θ | 1-3 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Bologni <i>et al.</i> (2021) | 2021 | CNN | C | S | Waveforms | θ, d | 1 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Cao <i>et al.</i> (2021) | 2021 | SA | R | S | Log-Mel spectrograms + intensity | x, y, z | 0-2 | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Castellini <i>et al.</i> (2021) | 2021 | MLP | R | S | real + imaginary CPS | x, y | 1-3 | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Diaz-Guerra <i>et al.</i> (2021b) | 2021 | CNN | R | S | SRP-PHAT power map | x, y, z | 1 | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Emmanuel <i>et al.</i> (2021) | 2021 | CNN + SA | R | S | Log-spectrograms + intensity | ACCDOA | 1 | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Gelderblom <i>et al.</i> (2021) | 2021 | MLP | C/R | S | GCC-PHAT | θ | 2 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Gonçalves Pinto <i>et al.</i> (2021) | 2021 | CNN | R | S | Magnitude CPS | x, y | 1-10 | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ |
| Grumiaux <i>et al.</i> (2021a) | 2021 | CRNN | C | S | Intensity | θ, ϕ | 1-3 | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Grumiaux <i>et al.</i> (2021b) | 2021 | CNN + SA | C | S | Intensity | θ, ϕ | 1-3 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Guirguis <i>et al.</i> (2020) | 2021 | TCN | R | S | Magnitude + phase spectrograms | x, y, z | 1 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Hammer <i>et al.</i> (2021) | 2021 | U-net | C | S | Phase map of the RTF between each mic pair | θ | ∞ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ |
| He <i>et al.</i> (2021a) | 2021 | Res. CNN | C | WS | Magnitude + phase spectrograms | θ | 1-4 | ✓/✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ |
| He <i>et al.</i> (2021b) | 2021 | CNN | R | S | Waveforms | x, y, z | 1 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Huang and Perez (2021) | 2021 | Res. CNN + SA | R | S | Waveforms | ACCDOA | 1 | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Komatsu <i>et al.</i> (2020) | 2021 | CRNN | R | S | FOA magnitude + phase spectrograms | θ, ϕ | 1 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Krause <i>et al.</i> (2020a) | 2021 | CNN | R | S | Magnitude + phase spectrograms | x, y, z | 1 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ |
| Krause <i>et al.</i> (2020b) | 2021 | CRNN | R | S | Misc. | θ, ϕ | 1 | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Lee <i>et al.</i> (2021a) | 2021 | U-Net | R | S | SRP power map | x, y | 1-3 | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| Lee <i>et al.</i> (2021b) | 2021 | CNN + attention | C | S | Log-Mel spectrograms + intensity | θ | 1 | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Liu <i>et al.</i> (2021) | 2021 | CNN | C | S | Intensity | θ | 1 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Naranjo-Alcazar <i>et al.</i> (2021) | 2021 | Res. CRNN | R | S | Log-Mel spectrograms + GCC-PHAT | ACCDOA | 1 | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Nguyen <i>et al.</i> (2021a) | 2021 | CRNN | C | S | Intensity/GCC-PHAT | θ, ϕ | 1 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Nguyen <i>et al.</i> (2021b) | 2021 | CNN + RNN/SA | R | S | Log-spectrograms + DRR + SCM eigenvectors | ACCDOA | 1 | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Park <i>et al.</i> (2021a) | 2021 | SA | R | S | log-Mel spectrograms + intensity | x, y, z | 1 | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Poschadel <i>et al.</i> (2021a) | 2021 | CRNN | C | S | HOA magnitude + phase spectrograms | θ, ϕ | 1 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Poschadel <i>et al.</i> (2021b) | 2021 | CRNN | C | S | HOA magnitude + phase spectrograms | θ, ϕ | 2-3 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ |
| Pujol <i>et al.</i> (2021) | 2021 | Res. CNN | R | S | Waveforms | θ, ϕ | 1 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✓ | ✓ |
| Rho <i>et al.</i> (2021) | 2021 | CRNN + SA | R | S | Log-Mel spectrograms + intensity | θ, ϕ | 1 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Schymura <i>et al.</i> (2021) | 2021 | CNN + SA | R | S | Magnitude + phase spectrograms | θ, ϕ | 1 | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ |
| Schymura <i>et al.</i> (2020) | 2021 | CNN + AE + attent. | R | S | FOA magnitude + phase spectrograms | θ, ϕ | 1 | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Shimada <i>et al.</i> (2021) | 2021 | Res. CRNN + SA | R | S | IPD | ACCDOA | 1 | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Subramanian <i>et al.</i> (2021a) | 2021 | CRNN | C/R | S | Phase spectrogram | θ | 2 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ |
| Subramanian <i>et al.</i> (2021b) | 2021 | CRNN | C | S | Phase spectrograms, IPD | θ | 2 | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ |
| Sudarsanam <i>et al.</i> (2021) | 2021 | SA | R | S | Log-Mel spectrograms + intensity | ACCDOA | 1 | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |

HOW

Prerak

JA

A survey
method

Pierre-Amar

¹Nantes Univer

²Orange Labs,

³Univ. Grenob

| Author | Year | Architecture | Type | Learn. | Input features | Output | Sources | | | Data | | | | | | | |
|--------------------------------------|------|--------------------|------|--------|---------------------------------------------|----------------|----------|------|------|-------|----|------|----|---|---|---|---|
| | | | | | | | NoS | Kno. | Mov. | Train | | Test | | | | | |
| | | | | | | | | | | SA | RA | SR | RR | | | | |
| Adavanne <i>et al.</i> (2021) | 2021 | CRNN + SA | R | S | FOA Mel spectrograms + intensity + GCC-PHAT | x, y, z | 2 | ✓ | ✓ | X | X | X | ✓ | X | X | X | ✓ |
| Bai <i>et al.</i> (2021) | 2021 | Res. CRNN | R | S | Log-Mel spectrograms + intensity | x, y, z | 1 | ✓ | ✓ | X | X | X | ✓ | X | X | X | ✓ |
| Bianco <i>et al.</i> (2021) | 2021 | VAE | C | SS | RTF | θ | 1 | ✓ | X | X | X | ✓ | X | X | X | ✓ | ✓ |
| Bohlender <i>et al.</i> (2021) | 2021 | CNN/CRNN | C | S | Phase map | θ | 1-3 | ✓ | X | X | X | ✓ | X | X | X | X | ✓ |
| Bologni <i>et al.</i> (2021) | 2021 | CNN | C | S | Waveforms | θ, d | 1 | ✓ | X | X | X | ✓ | X | X | X | ✓ | ✓ |
| Cao <i>et al.</i> (2021) | 2021 | SA | R | S | Log-Mel spectrograms + intensity | x, y, z | 0-2 | X | ✓ | X | X | X | ✓ | X | X | X | ✓ |
| Castellini <i>et al.</i> (2021) | 2021 | MLP | R | S | real + imaginary CPS | x, y | 1-3 | ✓ | X | ✓ | X | X | ✓ | X | X | X | ✓ |
| Diaz-Guerra <i>et al.</i> (2021b) | 2021 | CNN | R | S | SRP-PHAT power map | x, y, z | 1 | ✓ | ✓ | X | X | ✓ | X | X | X | ✓ | ✓ |
| Emmanuel <i>et al.</i> (2021) | 2021 | CNN + SA | R | S | Log-spectrograms + intensity | ACCDOA | 1 | ✓ | ✓ | X | X | X | ✓ | X | X | X | ✓ |
| Gelderblom <i>et al.</i> (2021) | 2021 | MLP | C/R | S | GCC-PHAT | θ | 2 | ✓ | X | X | X | ✓ | X | X | X | X | ✓ |
| Gonçalves Pinto <i>et al.</i> (2021) | 2021 | CNN | R | S | Magnitude CPS | x, y | 1-10 | X | X | ✓ | X | X | X | ✓ | X | X | X |
| Grumiaux <i>et al.</i> (2021a) | 2021 | CRNN | C | S | Intensity | θ, ϕ | 1-3 | ✓ | ✓ | X | X | ✓ | X | X | X | ✓ | ✓ |
| Grumiaux <i>et al.</i> (2021b) | 2021 | CNN + SA | C | S | Intensity | θ, ϕ | 1-3 | ✓ | X | X | X | ✓ | X | X | X | ✓ | ✓ |
| Guirguis <i>et al.</i> (2020) | 2021 | TCN | R | S | Magnitude + phase spectrograms | x, y, z | 1 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Hammer <i>et al.</i> (2021) | 2021 | U-net | C | S | Phase map of the RTF between each mic pair | θ | ∞ | X | ✓ | X | X | ✓ | X | X | X | X | ✓ |
| He <i>et al.</i> (2021a) | 2021 | Res. CNN | C | WS | Magnitude + phase spectrograms | θ | 1-4 | ✓/X | X | X | X | ✓ | ✓ | X | X | X | ✓ |
| He <i>et al.</i> (2021b) | 2021 | CNN | R | S | Waveforms | x, y, z | 1 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Huang and Perez (2021) | 2021 | Res. CNN + SA | R | S | Waveforms | ACCDOA | 1 | ✓ | ✓ | X | X | X | ✓ | X | X | X | ✓ |
| Komatsu <i>et al.</i> (2020) | 2021 | CRNN | R | S | FOA magnitude + phase spectrograms | θ, ϕ | 1 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Krause <i>et al.</i> (2020a) | 2021 | CNN | R | S | Magnitude + phase spectrograms | x, y, z | 1 | ✓ | X | X | X | ✓ | X | X | X | ✓ | X |
| Krause <i>et al.</i> (2020b) | 2021 | CRNN | R | S | Misc. | θ, ϕ | 1 | ✓ | X | X | X | ✓ | X | X | X | ✓ | X |
| Lee <i>et al.</i> (2021a) | 2021 | U-Net | R | S | SRP power map | x, y | 1-3 | X | X | ✓ | X | X | X | ✓ | X | X | ✓ |
| Lee <i>et al.</i> (2021b) | 2021 | CNN + attention | C | S | Log-Mel spectrograms + intensity | θ | 1 | ✓ | ✓ | X | X | ✓ | X | X | X | ✓ | ✓ |
| Liu <i>et al.</i> (2021) | 2021 | CNN | C | S | Intensity | θ | 1 | ✓ | X | X | X | ✓ | X | X | X | ✓ | ✓ |
| Naranjo-Alcazar <i>et al.</i> (2021) | 2021 | Res. CRNN | R | S | Log-Mel spectrograms + GCC-PHAT | ACCDOA | 1 | ✓ | ✓ | X | X | X | ✓ | X | X | X | ✓ |
| Nguyen <i>et al.</i> (2021a) | 2021 | CRNN | C | S | Intensity/GCC-PHAT | θ, ϕ | 1 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Nguyen <i>et al.</i> (2021b) | 2021 | CNN + RNN/SA | R | S | Log-spectrograms + DRR + SCM eigenvectors | ACCDOA | 1 | ✓ | ✓ | X | X | X | ✓ | X | X | X | ✓ |
| Park <i>et al.</i> (2021a) | 2021 | SA | R | S | log-Mel spectrograms + intensity | x, y, z | 1 | ✓ | ✓ | X | X | X | ✓ | X | X | X | ✓ |
| Poschadel <i>et al.</i> (2021a) | 2021 | CRNN | C | S | HOA magnitude + phase spectrograms | θ, ϕ | 1 | ✓ | X | X | X | ✓ | X | X | X | ✓ | ✓ |
| Poschadel <i>et al.</i> (2021b) | 2021 | CRNN | C | S | HOA magnitude + phase spectrograms | θ, ϕ | 2-3 | ✓ | X | X | X | ✓ | X | X | X | ✓ | ✓ |
| Pujol <i>et al.</i> (2021) | 2021 | Res. CNN | R | S | Waveforms | θ, ϕ | 1 | ✓ | X | X | X | ✓ | X | X | ✓ | ✓ | ✓ |
| Rho <i>et al.</i> (2021) | 2021 | CRNN + SA | R | S | Log-Mel spectrograms + intensity | θ, ϕ | 1 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Schymura <i>et al.</i> (2021) | 2021 | CNN + SA | R | S | Magnitude + phase spectrograms | θ, ϕ | 1 | ✓ | X | ✓ | X | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Schymura <i>et al.</i> (2020) | 2021 | CNN + AE + attent. | R | S | FOA magnitude + phase spectrograms | θ, ϕ | 1 | ✓ | X | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Shimada <i>et al.</i> (2021) | 2021 | Res. CRNN + SA | R | S | IPD | ACCDOA | 1 | ✓ | ✓ | X | X | X | ✓ | X | X | X | ✓ |
| Subramanian <i>et al.</i> (2021a) | 2021 | CRNN | C/R | S | Phase spectrogram | θ | 2 | ✓ | X | X | X | ✓ | X | X | X | ✓ | X |
| Subramanian <i>et al.</i> (2021b) | 2021 | CRNN | C | S | Phase spectrograms, IPD | θ | 2 | ✓ | X | X | X | ✓ | X | X | X | ✓ | X |
| Sudarsanam <i>et al.</i> (2021) | 2021 | SA | R | S | Log-Mel spectrograms + intensity | ACCDOA | 1 | ✓ | ✓ | X | X | X | ✓ | X | X | X | ✓ |

3 layers of **acoustic simulation realism** at train time

Wall Realism

Mic. Realism

Source Realism

3 layers of **acoustic simulation realism** at train time

Wall Realism

Mic. Realism

Source Realism

Naive: identical,
frequency-independent
walls

3 layers of **acoustic simulation realism** at train time

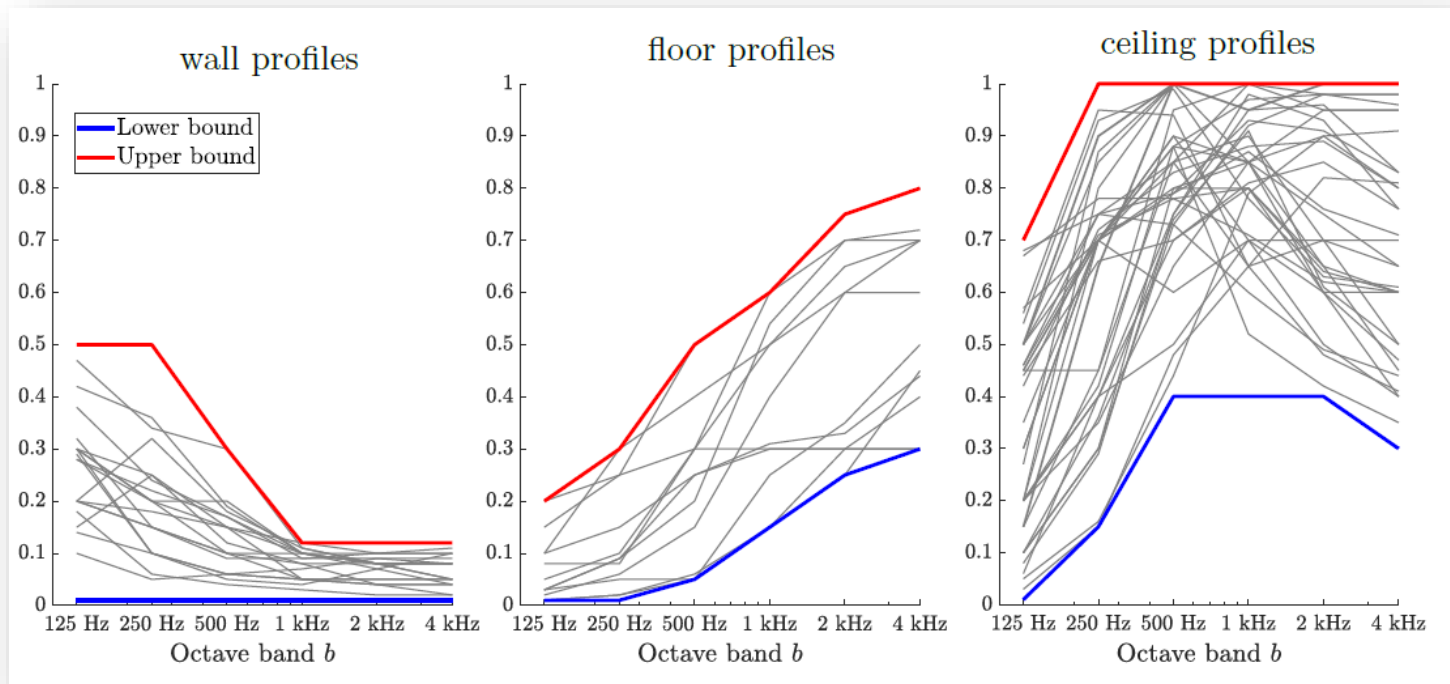
Wall Realism

Mic. Realism

Source Realism

Naive: identical, frequency-independent walls

Advanced: Based on real material absorption databases



3 layers of **acoustic simulation realism** at train time

Wall Realism

Mic. Realism

Source Realism

Naive: identical,
frequency-independent
walls

Naive: omnidirectional, frequency-
independent microphones and sources

Advanced:
Based on real
material
absorption
databases

A new preprint

3 layers of **acoustic simulation realism** at train time

Wall Realism

Naive: identical, frequency-independent walls

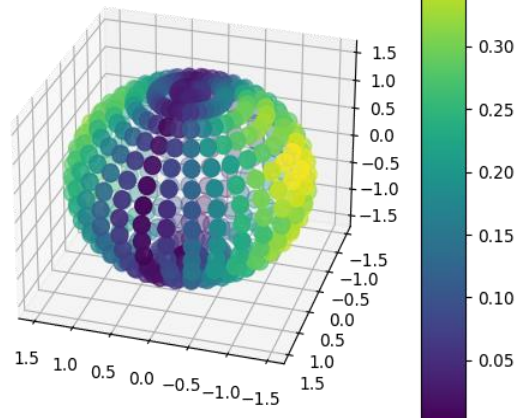
Advanced: Based on real material absorption databases

Mic. Realism

Naive: omnidirectional, frequency-independent microphones and sources

Advanced: Based on real measured directivity profiles

Figure of Eight Receiver , Mag|STFT|
@ 2 kHz



3 Real Test Sets

DIRHA (2017)

VoiceHome2 (2018)

STARSS (2022)

- > 20 real Human speakers
- > 15 real rooms (living room, kitchen, classroom)
- 3 arrays of 2 microphones

95 minutes of speech annotated with sound source direction

Results

- Advanced simulation significantly improves sound source localization results across all test sets
- Every added layer of realism contributes to these results

| Real Test Sets | VoiceHome-2 [24] | | DIRHA [25] | | STARS22 [26] | |
|-----------------------|------------------|------------------|------------|------------------|--------------|-------------------|
| Methods | ↑ Recall | ↓ MAE (°) | ↑ Recall | ↓ MAE (°) | ↑ Recall | ↓ MAE (°) |
| SRP-PHAT | 70% | 9.9 ± 1.5 | 61% | 15.0 ± 2.3 | 45% | 14.9 ± 0.6 |
| Naive Training | 78% | 7.6 ± 1.2 | 77% | 8.4 ± 1.4 | 57% | 12.9 ± 0.6 |
| Advanced Training | 85% | 5.8 ± 0.8 | 84% | 6.3 ± 1.0 | 61% | 11.4 ± 0.5 |
| Ablation study | | | | | | |
| w/o wall realism | 83% | 6.2 ± 0.8 | 81% | 7.5 ± 1.4 | 59% | 12.1 ± 0.6 |
| w/o source realism | 82% | 7.1 ± 1.1 | 80% | 7.8 ± 1.2 | 63% | 11.4 ± 0.6 |
| w/o receiver realism | N/A | N/A | 78% | 8.3 ± 1.5 | 53% | 13.4 ± 0.6 |

UNIVERSITÉ DE GRENOBLE

THÈSE

Pour obtenir le grade de

DOCTEUR DE L'UNIVERSITÉ DE GRENOBLE

Spécialité : **Mathématiques et Informatique**

Arrêté ministériel :

Présentée par

Antoine Deleforge

Thèse dirigée par **Radu Horaud**

préparée au sein de l'Université Joseph Fourier, de l'INRIA Grenoble Rhône-Alpes et de L'École Doctorale de Mathématiques, Sciences et Technologies de l'Information, Informatique

Acoustic Space Mapping

A Machine Learning Approach to Sound Source Separation and Localization

Thèse soutenue publiquement le ,
devant le jury composé de :

Pr. Jonathon Chambers

Loughborough University, Rapporteur

Pr. Rémi Gribonval

INRIA Rennes, Rapporteur

Pr. Florence Forbes

INRIA Grenoble Rhône-Alpes, Examinatrice

Pr. Geoff MacLachlan

University of Queensland, Examineur

Pr. Laurent Girin

GIPSA-lab, Grenoble, Examineur

Pr. Radu Horaud

INRIA Grenoble Rhône-Alpes, Directeur de thèse



is possible. This could include live music recording on a stage or concert hall, speech localization, diarisation or enhancement in a meeting or conference room, or hearing aid devices (the system could be calibrated for a specific wearer).

6.2 Direction for Future Research

Rather than an end, we like to view this thesis as a starting point for fascinating future research topics. We propose here a non-exhaustive list of possible follow-ups.

- An important direction is to study more thoroughly the influence of changes in experimental conditions on binaural manifolds. What happens when changing the position of the recording setup? Moving to another room? What is the influence of the sound source distance and directivity? What happens when the HRTF change? While the PLOM model is a possible direction to improve robustness to such situations, other methods such as *transfer learning* [Pan 10] could be envisioned. A more ambitious idea would be to learn acoustic spaces in virtual environments, using a room simulator such as Roomsim [Campbell 05]. One could imagine learning many different models in different room configurations, *e.g.*, microphones position, room size, reverberations. When dealing with real world data, the most appropriate model could be selected from virtually learned one using, *e.g.*, model selection.
- In our view, the surprisingly good results obtained so far could help to open the doors to a new category of binaural processing algorithms for a deeper understanding. First of all...

THÈSE

Pour obtenir le grade de

DOCTEUR DE L'UNIVERSITÉ

Spécialité : **Mathématiques**

Arrêté ministériel :

Présentée par

Antoine Deleforge

Thèse dirigée par Raouf

préparée au sein de l'Université de Grenoble Alpes et de L'École Doctorale de l'Information, des Télécommunications et de la Signalétique

Acoustic Source Separation

A Machine Learning Approach

Thèse soutenue publiquement devant le jury composé de

Pr. Jonathon Chan
Loughborough University

Pr. Rémi Gribonval
INRIA Rennes, Rapporteur

Pr. Florence Forbes
INRIA Grenoble Rhône-Alpes

Pr. Geoff MacLachlan
University of Queensland

Pr. Laurent Girin
GIPSA-lab, Grenoble, Président

Pr. Radu Horaud
INRIA Grenoble Rhône-Alpes



What's next?

2023 +

What's next?

2023 +

What's next?

2023 +



What's next?

2023 +





« *What is the shape of
the room?* »



« *What is the shape of the room?* »

« *Is the floor made of tiles or carpet?* »



« *What is the shape of the room?* »

« *Is the floor made of tiles or carpet?* »

- Room acoustic diagnosis
- Audio augmented reality
- Echo-aware audio signal enhancement



« *What is the shape of the room?* »

« *Is the floor made of tiles or carpet?* »

- Room acoustic diagnosis
- Audio augmented reality
- Echo-aware audio signal enhancement
- Plenty of interesting open inverse problems

Thank You Radu

for teaching me
the joys of research!

Thank You Radu

for teaching me
the joys of research!

