# Machine learning and applications to ocean acoustics

Peter Gerstoft,

**http://noiselab.ucsd.edu/**.  Slides and 42-page review paper [Bianco 2019]
With help from Mike Bianco, Emma Ozanich, Haiqiang Niu, Kay Gemba,
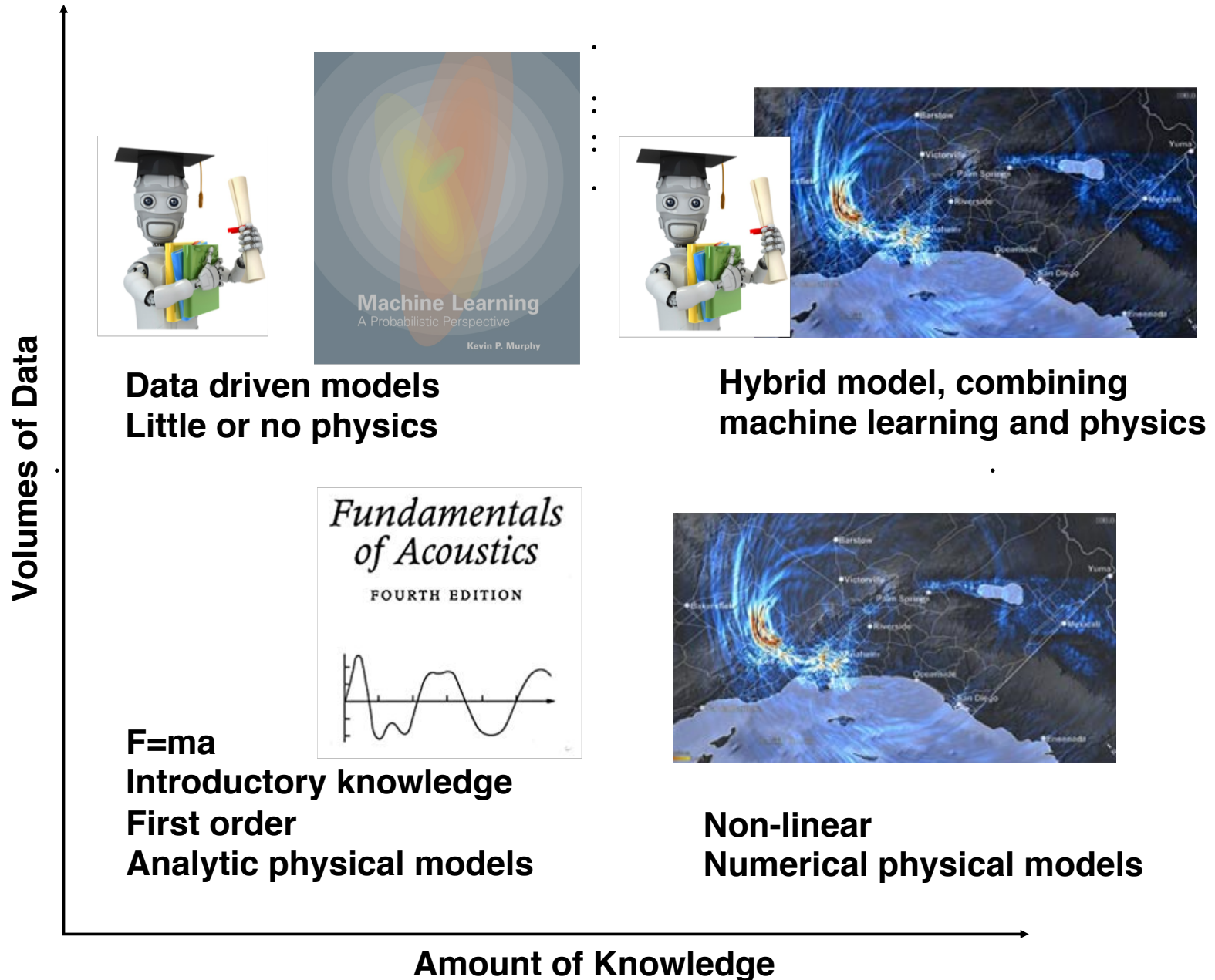James Traer, Christoph Mecklenbrauker, Eliza Michalopoulou

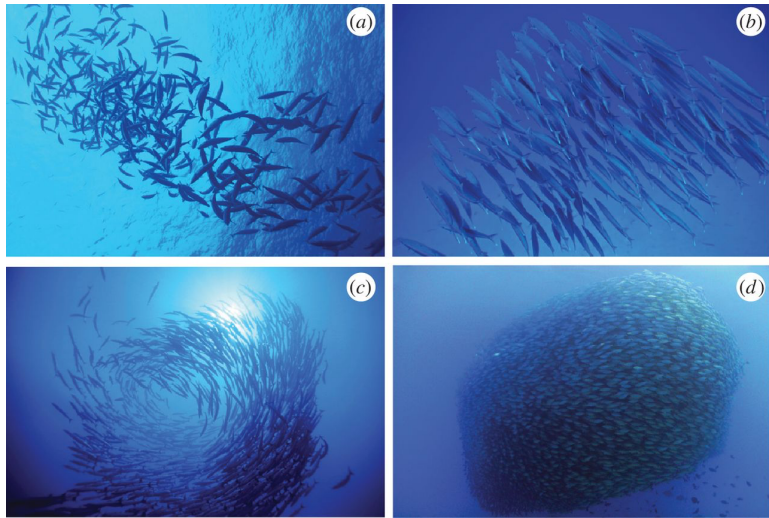Machine leaning  contains the mathematical tools we need to do
**data science**

Can Machine Learning
- Replace CTBTO/SONAR processing chain?
- Discover PDE (Partial differential equation) in video?
- Find sea mines?
- Design metamaterials?
- Predict earthquakes?
- Source location in the ocean waveguide w/o training?
- Replace 50 years of array processing?
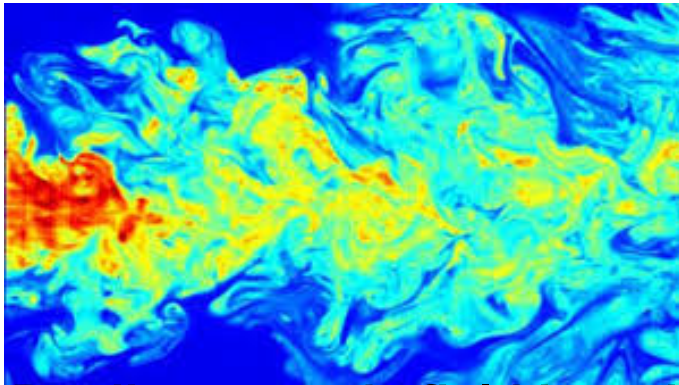- Learn the physical model (sound speed, temperature…)

# Machine learning versus knowledge based

Acoustic insight can be improved by leveraging the strengths of both physical and ML-based, data-driven models.
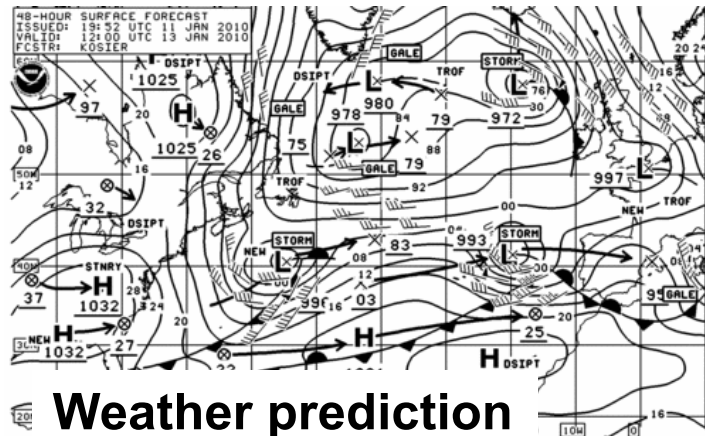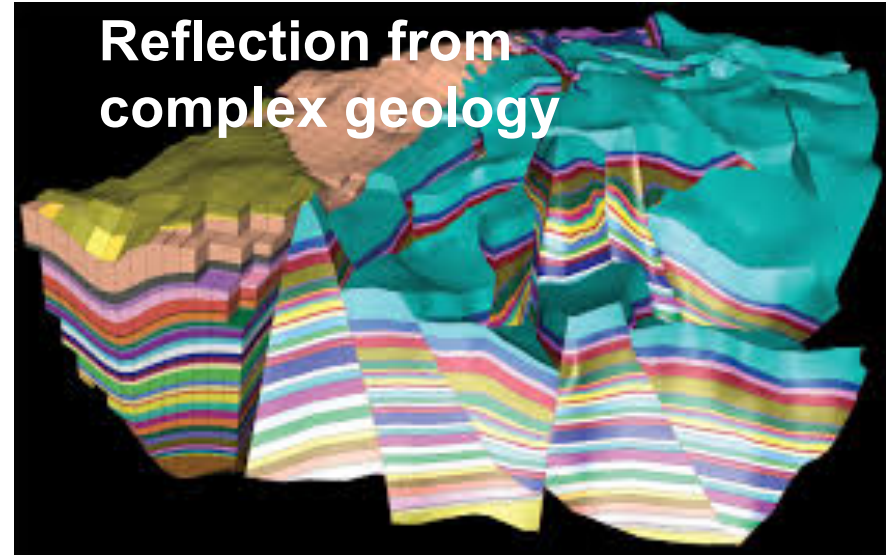


**Volumes of Data** (y-axis)

**Amount of Knowledge** (x-axis)

**Data driven models**
**Little or no physics**

**Hybrid model, combining**
**machine learning and physics**

**F=ma**
**Introductory knowledge**
**First order**
**Analytic physical models**

**Non-linear**
**Numerical physical models**

# We can't model everything...

Back scattering from fish school

Predict acoustic field in turbulence

Weather prediction

Reflection from complex geology

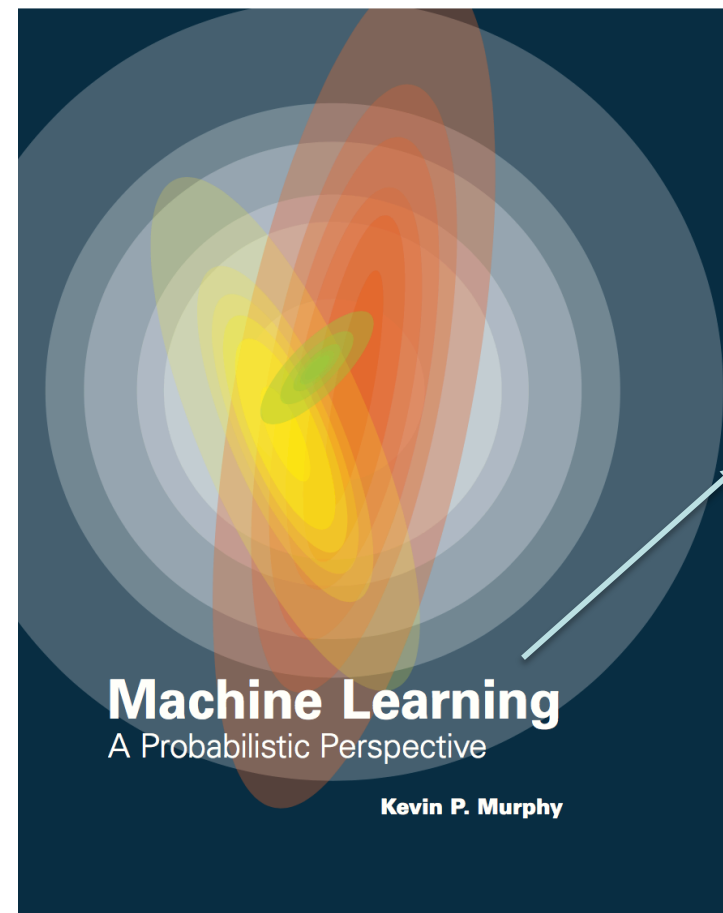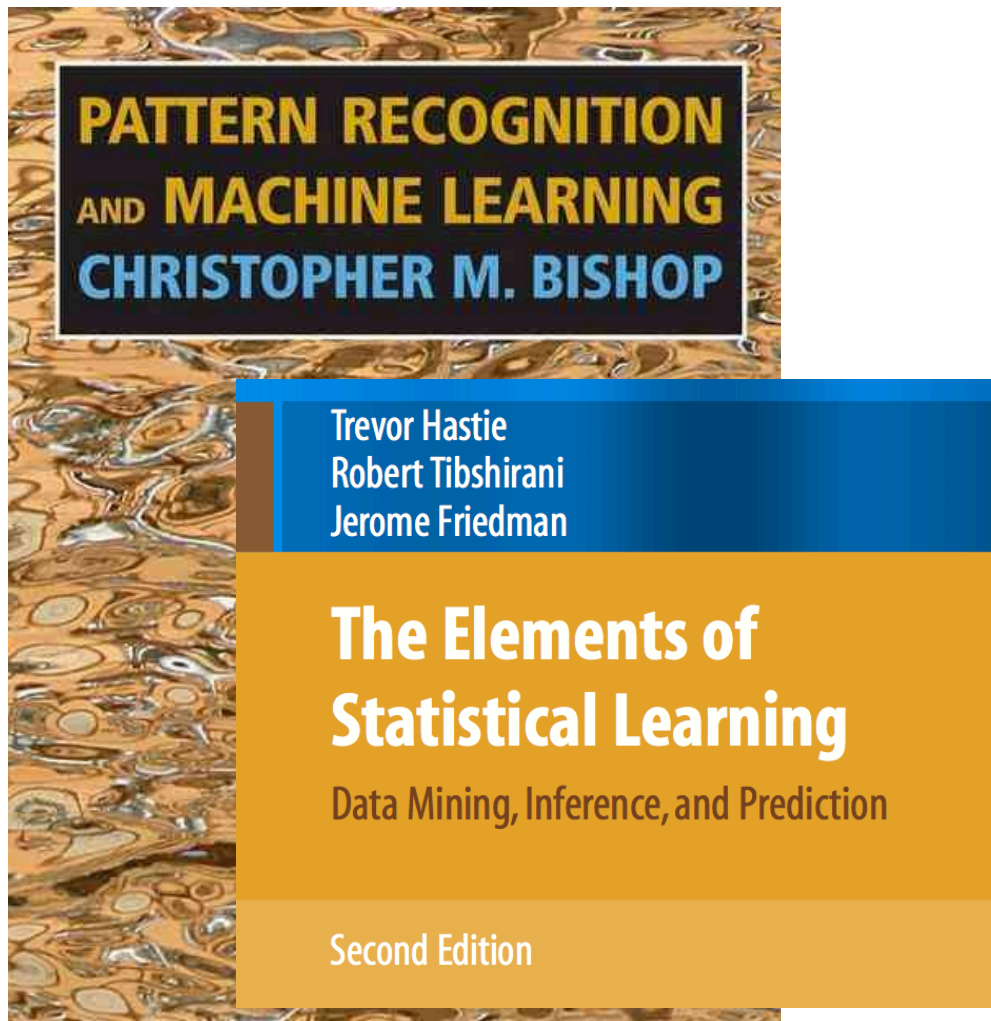Detection of mines. Navy uses dolphins to assist in this.

Dolphins = real ML!

# Machine Learning for physical Applications
## noiselab.ucsd.edu

Murphy: "…**the best way to make machines that can learn from data is to use the *tools of probability theory*, which has been the mainstay of statistics and engineering for centuries.**"
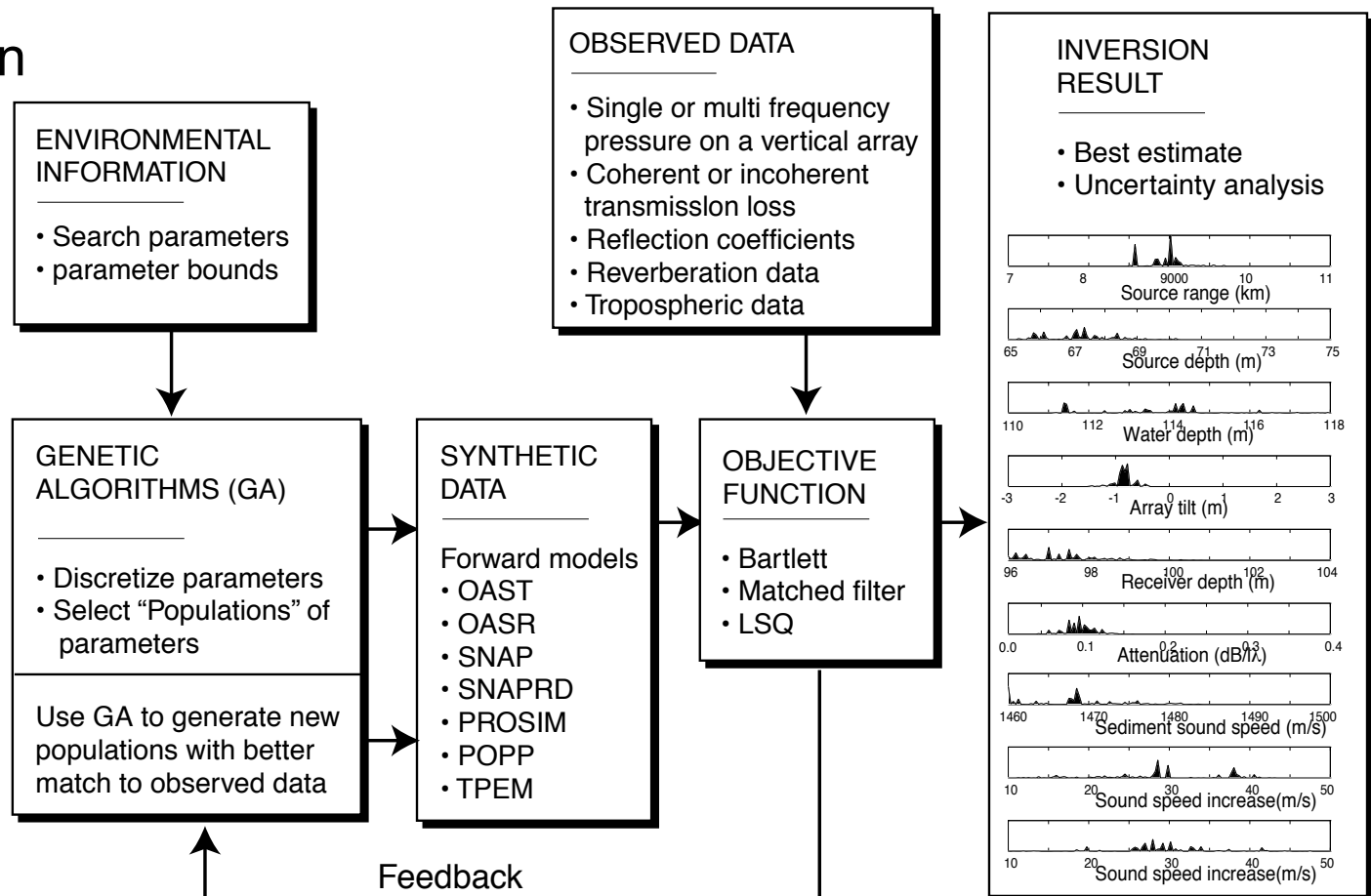


4

# SAGA (NURC 1992-97) is also ML

SAGA has the features that characterize a ML approach:

- Data-driven
- Model based
- Gaussian based likelihoods.
- Bayesian posterior probabilities

Also later additions

- Sequential estimation
- Particle Filtering

Gerstoft, 1994

**ENVIRONMENTAL INFORMATION**
- Search parameters
- parameter bounds

**OBSERVED DATA**
- Single or multi frequency pressure on a vertical array
- Coherent or incoherent transmisslon loss
- Reflection coefficients
- Reverberation data
- Tropospheric data

**INVERSION RESULT**
- Best estimate
- Uncertainty analysis

**GENETIC ALGORITHMS (GA)**
- Discretize parameters
- Select "Populations" of parameters

Use GA to generate new populations with better match to observed data

**SYNTHETIC DATA**

Forward models
- OAST
- OASR
- SNAP
- SNAPRD
- PROSIM
- POPP
- TPEM

**OBJECTIVE FUNCTION**
- Bartlett
- Matched filter
- LSQ

Feedback

Source range (km)
Source depth (m)
Water depth (m)
Array tilt (m)
Receiver depth (m)
Attenuation (dB/λ)
Sediment sound speed (m/s)
Sound speed increase(m/s)
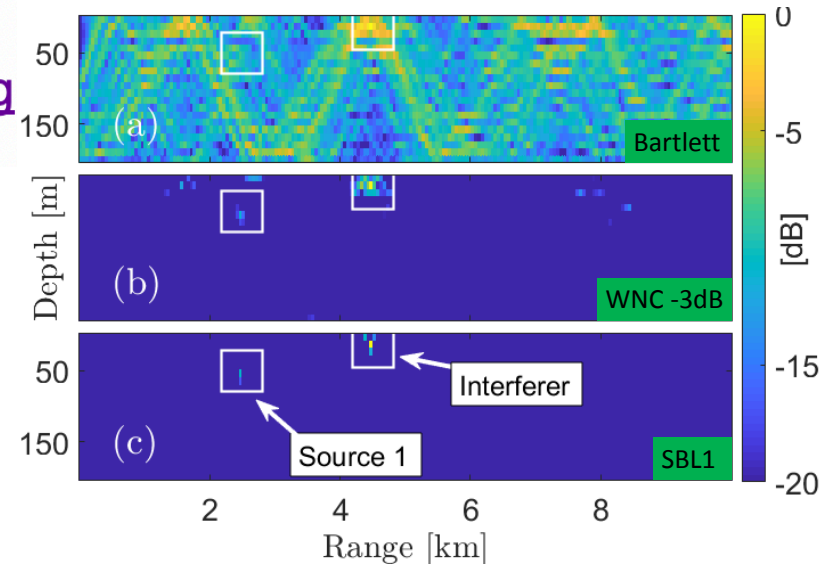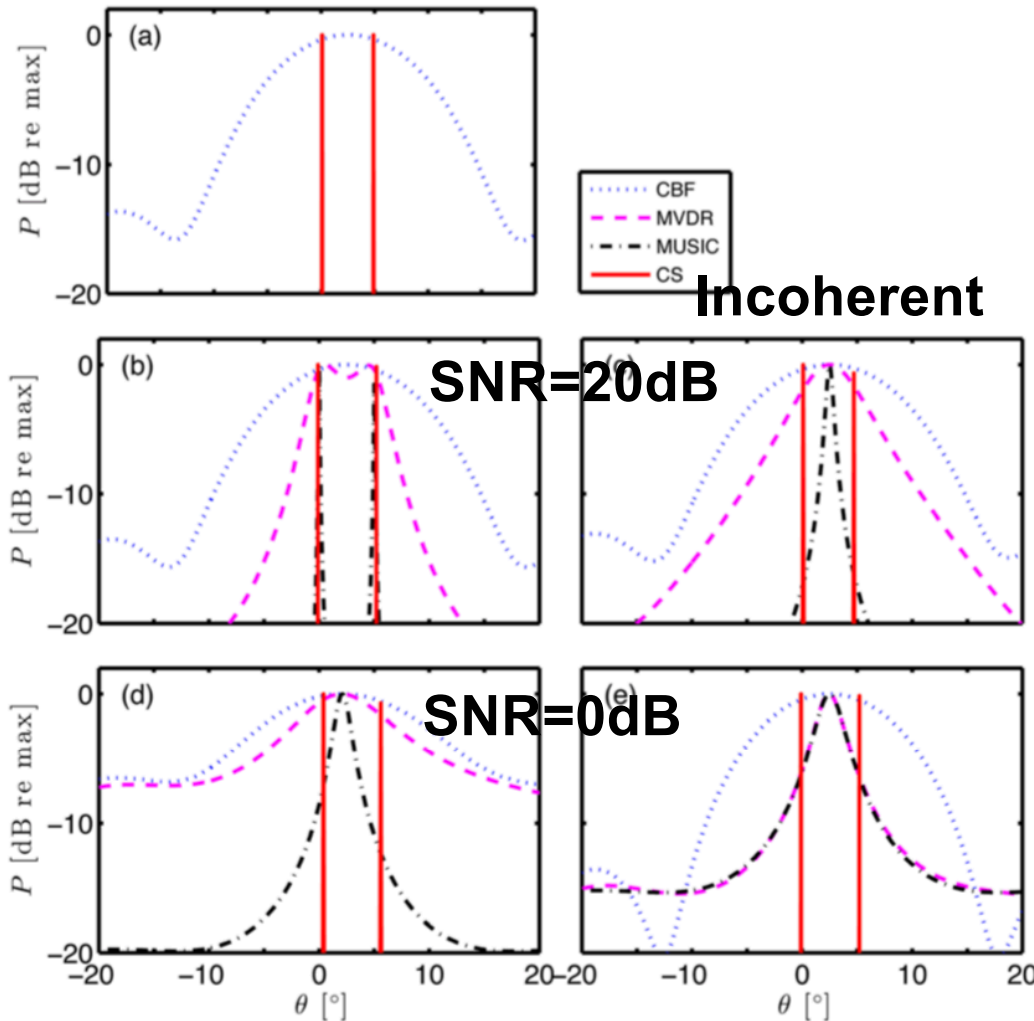Sound speed increase(m/s)

# Compressive beamforming is also ML

Compressive beamforming
A Xenaki, P Gerstoft, K Mosegaard - JASA, 2014 Cited by 142
Multiple and single snapshot compressive beamforming
P Gerstoft, A Xenaki, CF Mecklenbrauker- JASA, 2015 Cited by 78

**Coherent**

**Incoherent**

**SNR=20dB**

**SNR=0dB**

CBF
MVDR
MUSIC
CS

CS beamforming:
- **single or multiple snapshots**
- **coherent or incoherent**

Xenaki 2014, 2015, Gemba 2017

# Machine learning in acoustics: a review

Michael J. Bianco,[1, a)] Peter Gerstoft,[1] James Traer,[2] Emma Ozanich,[1] Marie A. Roch,[3] Sharon Gannot
Charles-Alban Deledalle,[5] and Weichang Li[6]

[1] *Scripps Institution of Oceanography, University of California San Diego, La Jolla, CA 92037, USA*

[2] *Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA 02139*

[3] *Department of Computer Science, San Diego State University, San Diego, CA 92182, USA*

[4] *Faculty of Engineering, Bar-Ilan University, Ramat-Gan 5290002, Israel*

[5] *Department of Electrical and Computer Engineering, University of California San Diego, La Jolla, CA 92037, USA*

[6] *Aramco Research Center-Houston, Aramco Services Company, Houston, TX 77084*

- **42-page** JASA review of ML theory. Available on **arXiv** or **http://noiselab.ucsd.edu/**.  **(**Pdf of talk is also there)

- Sections:
  - Machine learning principles
    - Supervised/ Unsupervised learning
  - Deep learning
  - Source localization in speech processing
  - Source localization in ocean acoustics
  - Bioacoustics
  - Seismic exploration
  - Perception of everyday sounds
    - Reverberation
    - Environmental sounds

# ML Principles

In ML, we are often interested in training a model to produce a desired output given inputs,

$$y=f(x) + \epsilon$$

- Input $\mathbf{x} \in \mathbb{R}^N$, N features

- output $\mathbf{y} \in \mathbb{R}^P$, P outputs

- **Supervised learning**: the P outputs have labelled examples (response variables $\mathbf{y}$)

- **Unsupervised learning**: there are no labels. The goal is to find interesting properties from $\mathbf{x}$, as an autoencoder $\tilde{\mathbf{x}}=\mathbf{f(x)}$

- ….. and we *train* the model



Features          Output

# Two ways to make computers do what you want:

**In Image processing this has been done:**

1) Hand-engineered design: Consciously figure out exactly how to manipulate symbolic representations to perform the task and then tell the computer in detail what to do.



Find edges

2) **Learning:** Show computers lots of examples of input with desired outputs. Let the computer learn how to map inputs to outputs using **general purpose** learning procedure
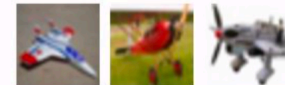


A close-up of a child holding a stuffed animal.

Input is an image

Output is a caption

**Example training set**
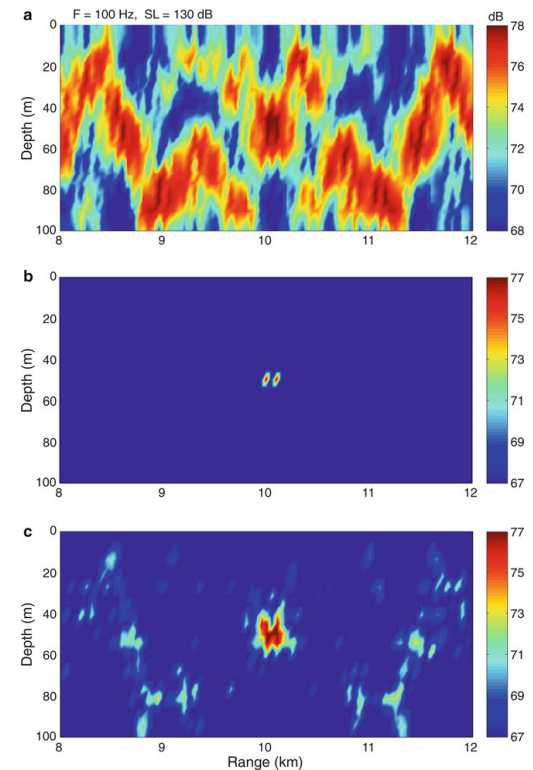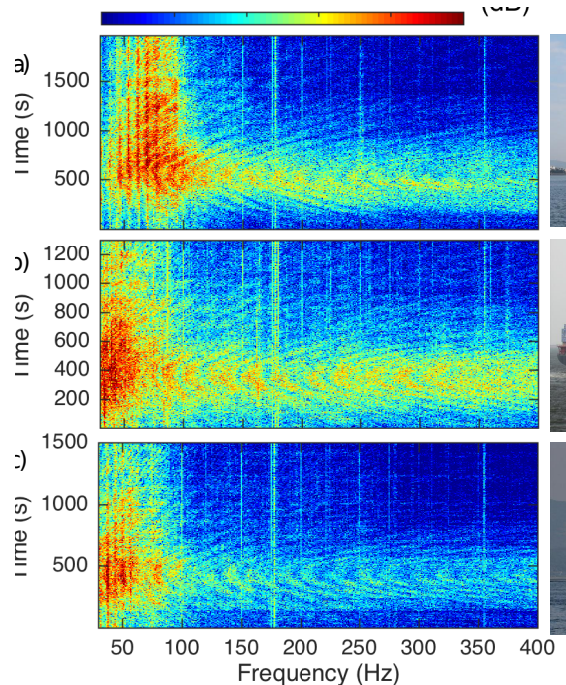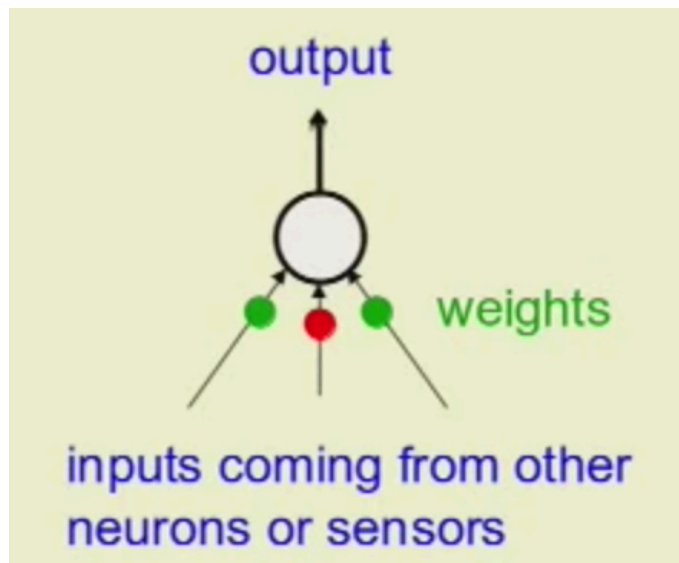


airplane
automobile
bird
cat
deer

# Two ways to make computers do what you want:

**In Ocean acoustics:**

1) Hand-engineered design: See the **1000 papers** on Match Field Processing! Sometimes it works…

=> **Old School**

2) Learning:   Show computers lots of examples of input with desired outputs. Let the computer learn how to map inputs to outputs using **general purpose** learning procedure
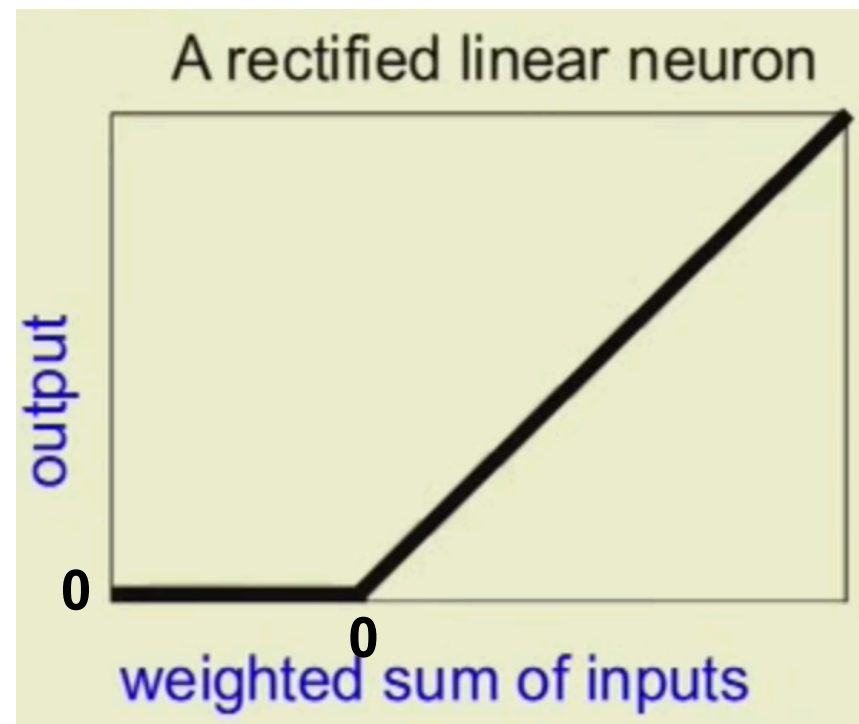
# What is an artificial neuron?

We simplify a real neuron to investigate how neurons can do computations that are too difficult to program as
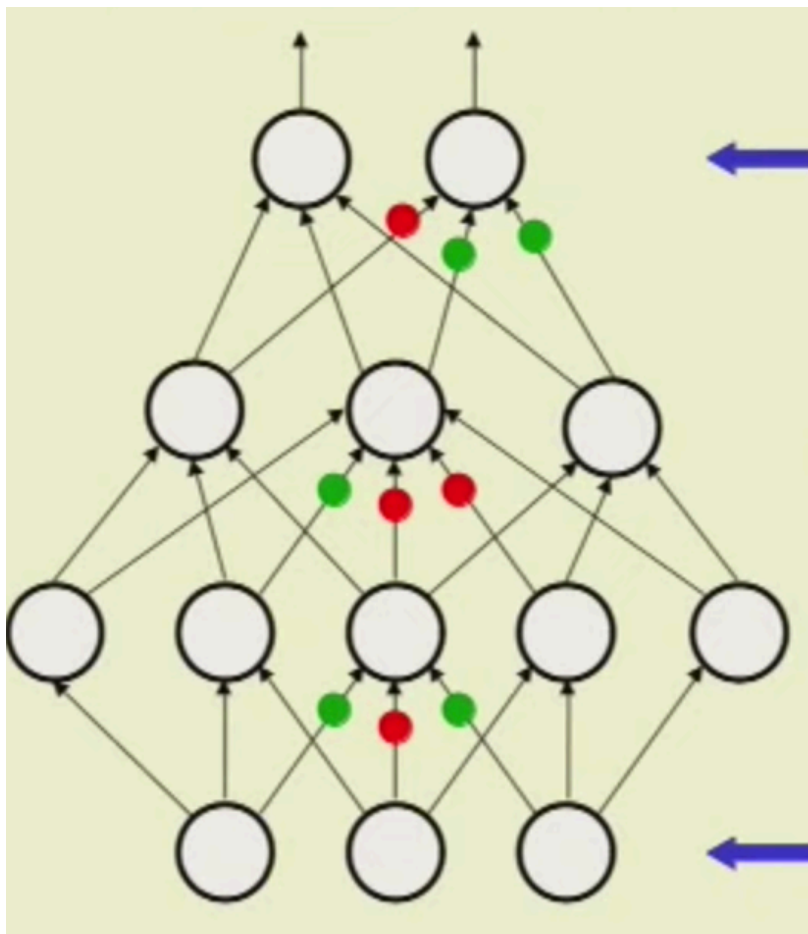- Converting image pixel intensity into string of words describing it

output

weights

inputs coming from other neurons or sensors

## ReLu

A rectified linear neuron

output

0

0

weighted sum of inputs

# What is artificial neural network

Connecting neurons in layers with no cycles gives a feed-forward neural net (FNN).

$$a_j = \text{ReLu}(\boldsymbol{w}^T \boldsymbol{x}) = \text{ReLu}(\sum_{n=1}^{N} w_n x_n)$$
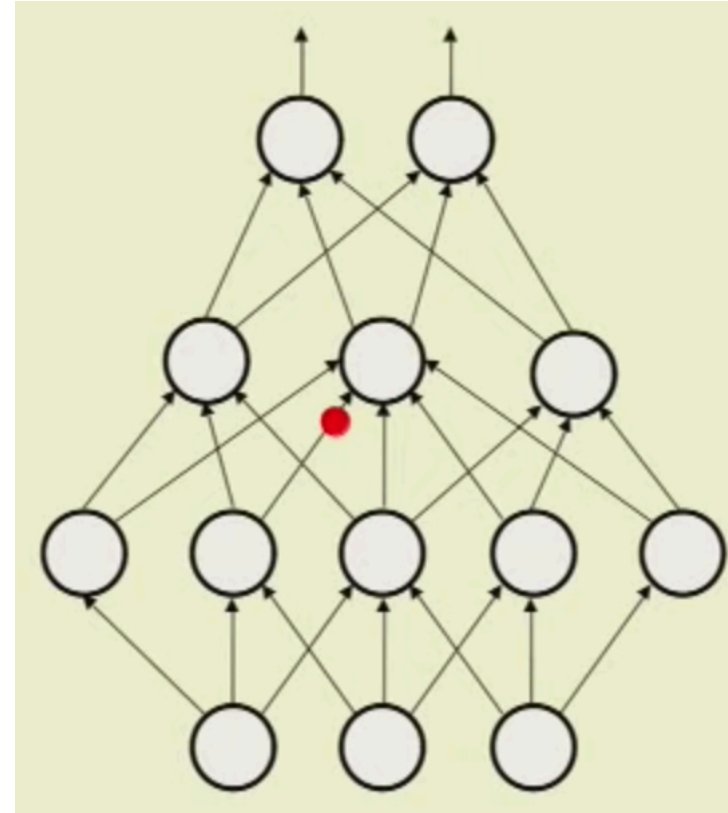


Output neurons

Multiple layers of hidden neurons
Hidden layers

Input neurons

# Supervised training vs backpropagatoin

Supervised training is inefficient:
- Take a few of the training cases and measure the NN output. (called **stochastic sampling**)
- Change one weight slightly.
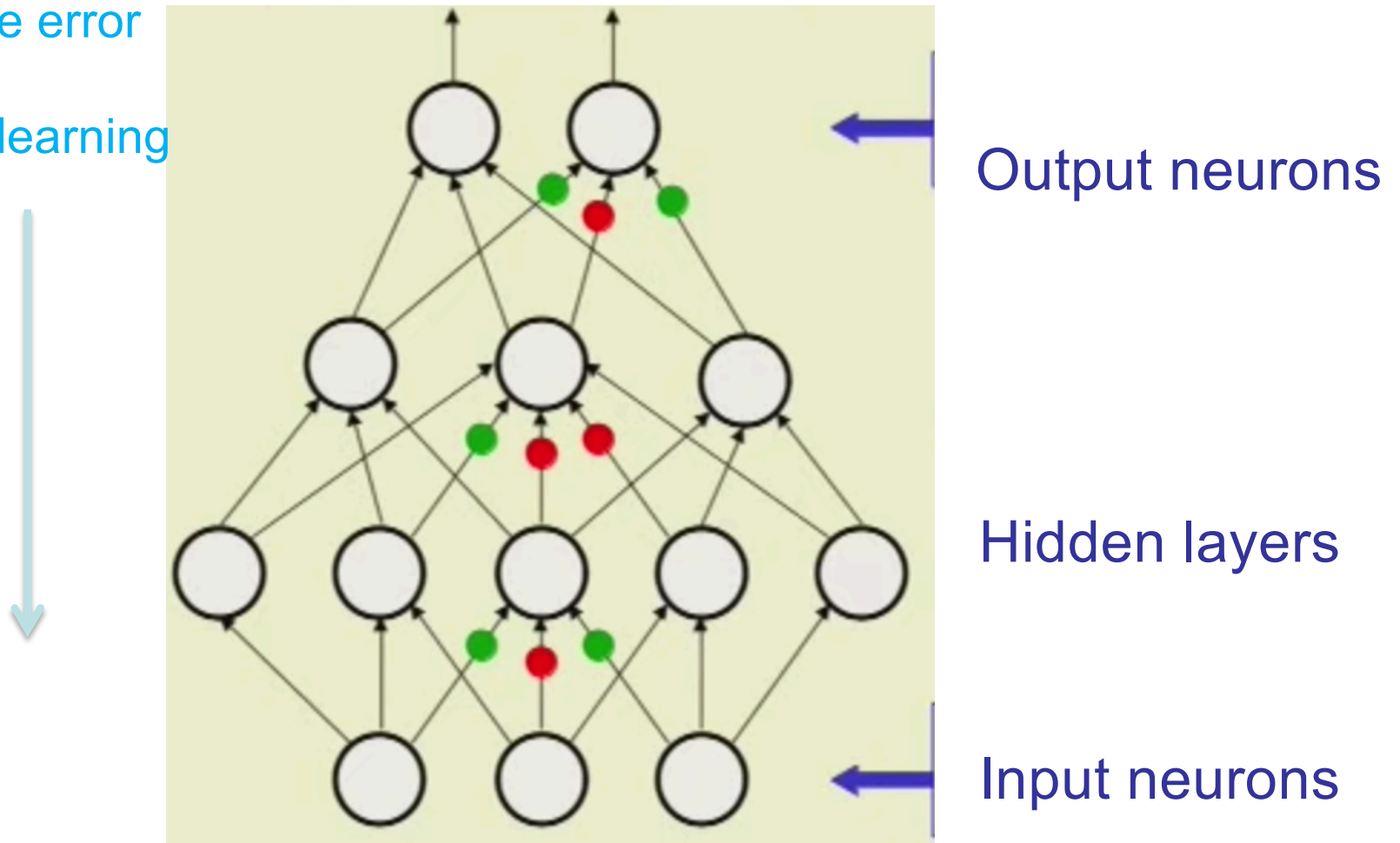- If NN output improved, **keep it.**



- **Backpropagation** efficiently compute how a change in weight effects the NN output.

- The error gradients for all of the weights is obtained at once. The chain rule dictates how the NN output change for each weight.

# How to learn many layers of features

Compare outputs with the correct answer to get the error signal

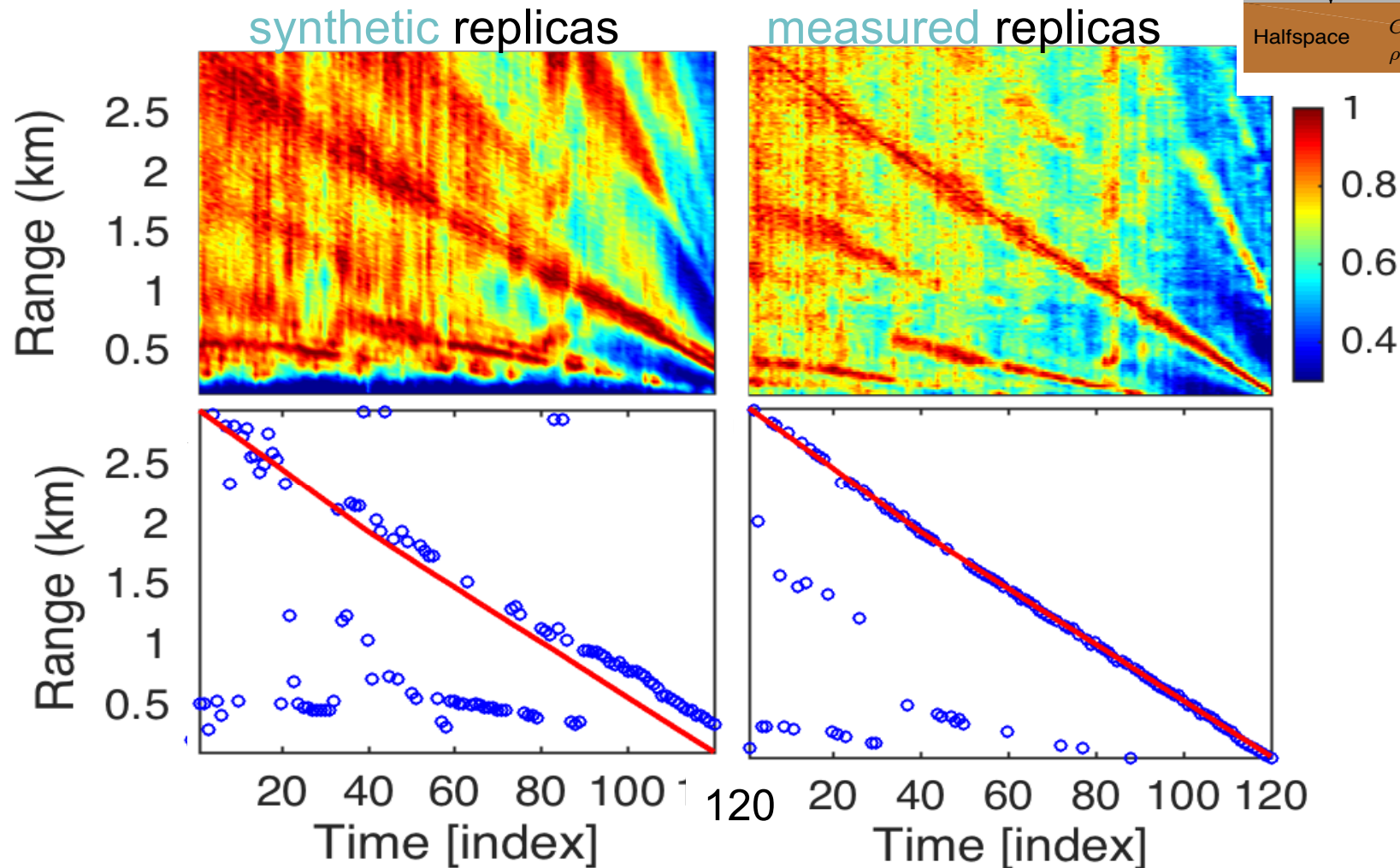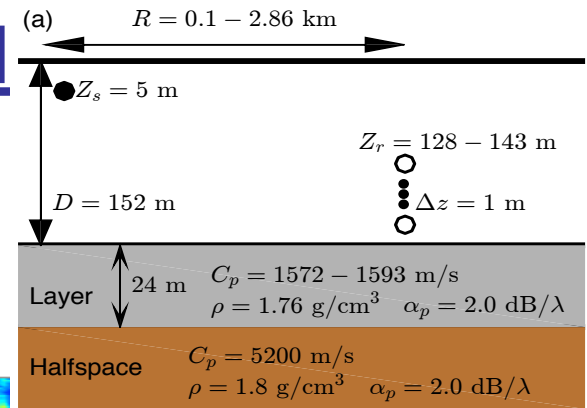Back-propagate error signal to get derivatives for learning



Output neurons

Hidden layers

Input neurons

For $L$ layers with $N$ neurons, we have $N^2L$ weights

# Matched-Field Processing on test data 1

Noise09 Frequencies [300:10:950]Hz

$$B = \mathbf{p}^H \mathbf{Cp}$$

$$E_{\text{MAPE}} = \frac{100}{N} \sum_{i=1}^{N} \left| \frac{Rp_i - Rg_i}{Rg_i} \right|$$

(a) $R = 0.1 - 2.86$ km

$Z_s = 5$ m

$Z_r = 128 - 143$ m

$D = 152$ m

$\Delta z = 1$ m

Layer — 24 m — $C_p = 1572 - 1593$ m/s  $\rho = 1.76$ g/cm$^3$  $\alpha_p = 2.0$ dB/$\lambda$

Halfspace  $C_p = 5200$ m/s  $\rho = 1.8$ g/cm$^3$  $\alpha_p = 2.0$ dB/$\lambda$

synthetic replicas

measured replicas



Mean Absolute Percentage Error error of MFPs:  **55%** and **19%**

Niu 2017a, JASA

# Pressure data preprocessing

Sound pressure

$$\mathbf{p}(f) = S(f)\mathbf{g}(f,\mathbf{r}) + \mathbf{n},$$

$S(f)$  Source term

Normalize pressure to reduce the effect of $|S(f)|$

$$\tilde{\mathbf{p}}(f) = \frac{\mathbf{p}(f)}{\sqrt{\sum_{l=1}^{L}|p_l(f)|^2}} = \frac{\mathbf{p}(f)}{\|\mathbf{p}(f)\|_2}$$
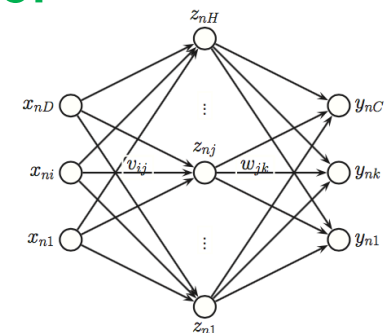
$L$  Number of sensors

Sample Covariance Matrix to reduce effect of source phase

$$\mathbf{C}(f) = \frac{1}{N_s}\sum_{s=1}^{N_s}\tilde{\mathbf{p}}_s(f)\tilde{\mathbf{p}}_s^H(f)$$

$N_s$  Number of snapshots

SCM is a conjugate symmetric matrix.

**Input vector X to NN: the real and imaginary parts of the entries of diagonal and upper triangular matrix in** $\mathbf{C}(f)$

# ML source range classification

Array Data: 300–950Hz with 10Hz increment, i.e., 66 frequencies.
16 hydrophones with 1 m spacing

First NN is trained with one source

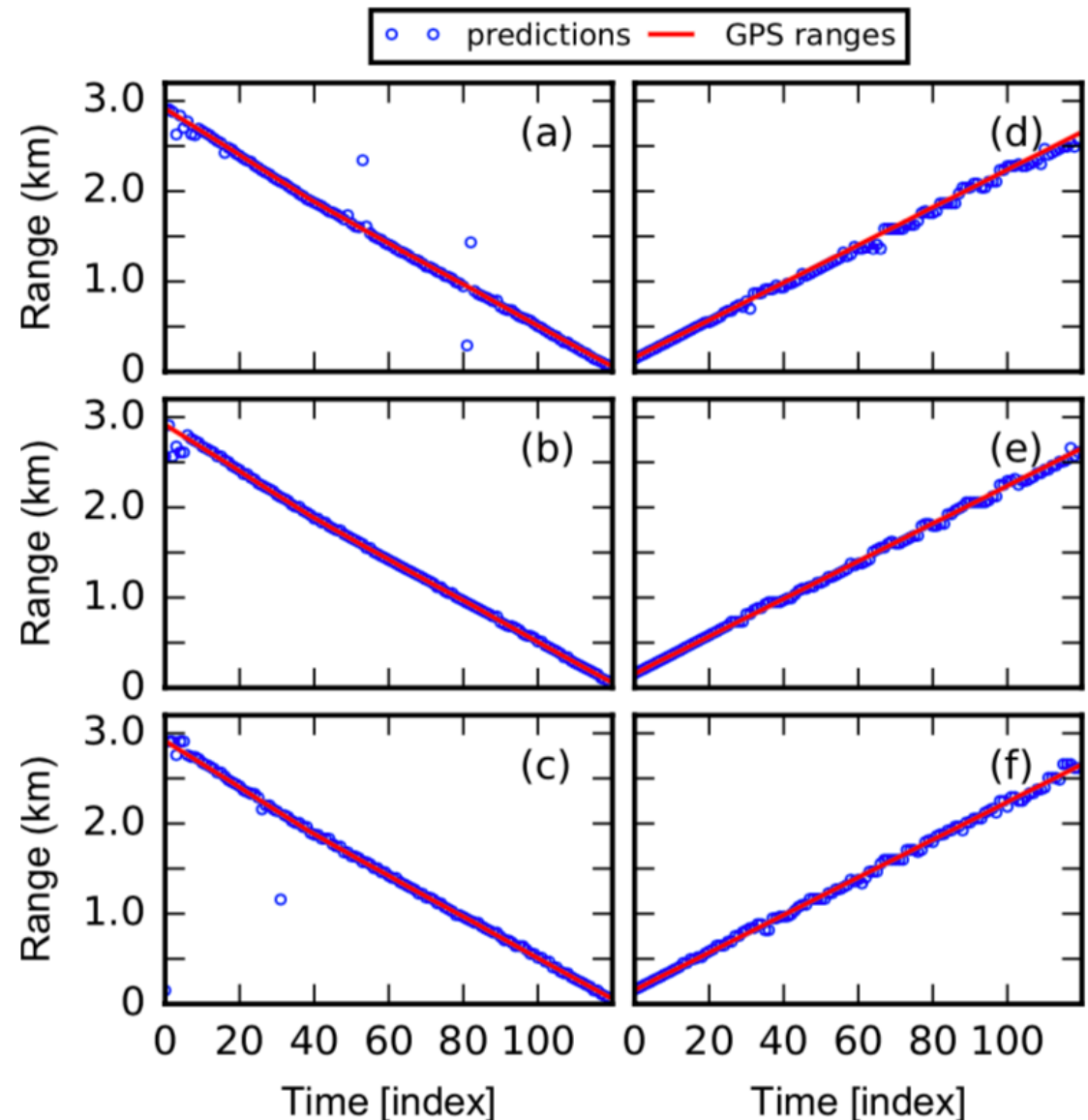**Test-Data-1**    **Test-Data-2**

**FNN**

3 hidden layers with 512 nodes

**SVM**

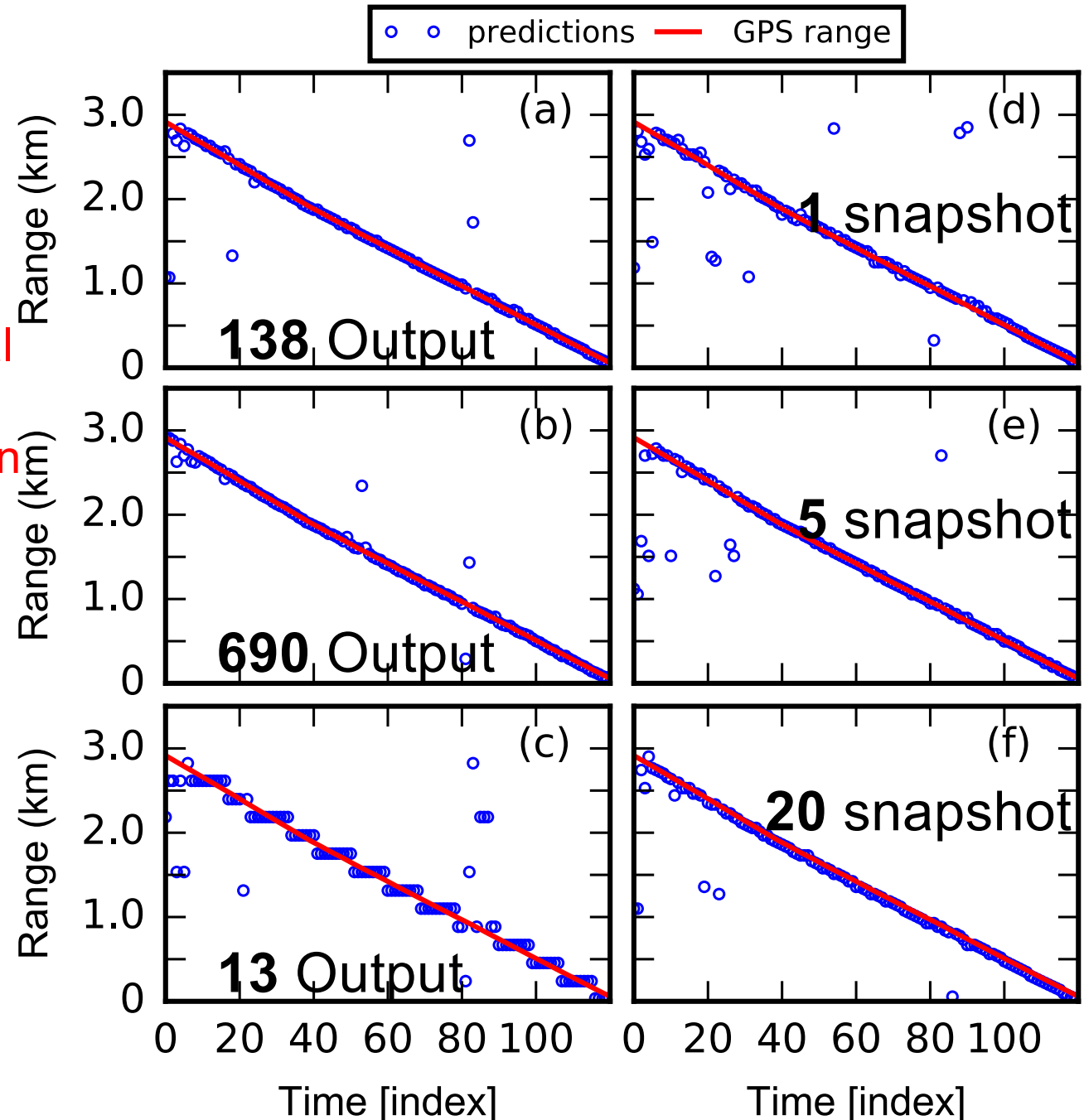Radial basis Functions

**RF**
Random forest



Niu 2017a, JASA

# Other parameters: FNN for range classification



**Conclusion**
- Easier than conventional MFP
- Classification easier than regression
- **FNN, SVM, RF** works.
- Works for:
  - multiple ships,
  - Deep/shallow water

60s Science
Scientfic Am

Niu 2017a, JASA

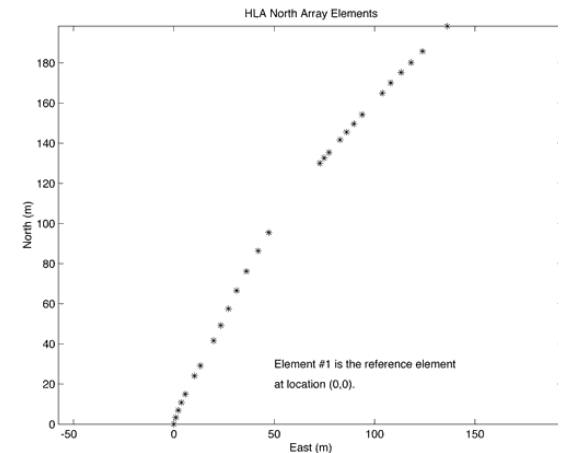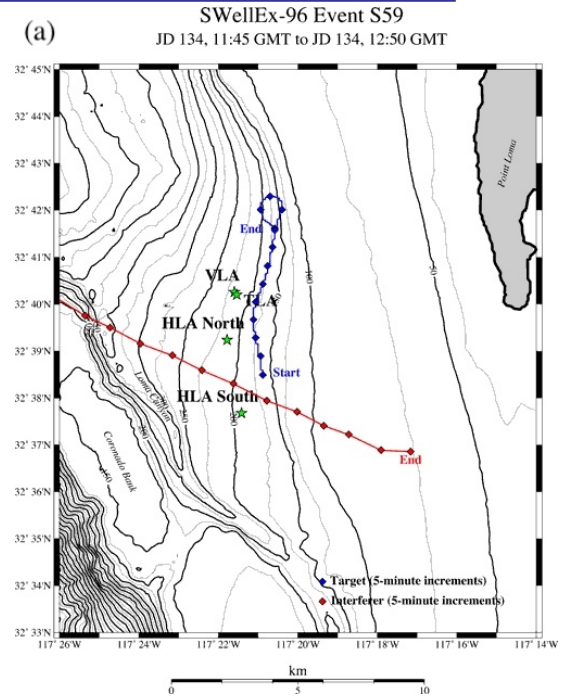# DOA estimation with Neural Networks
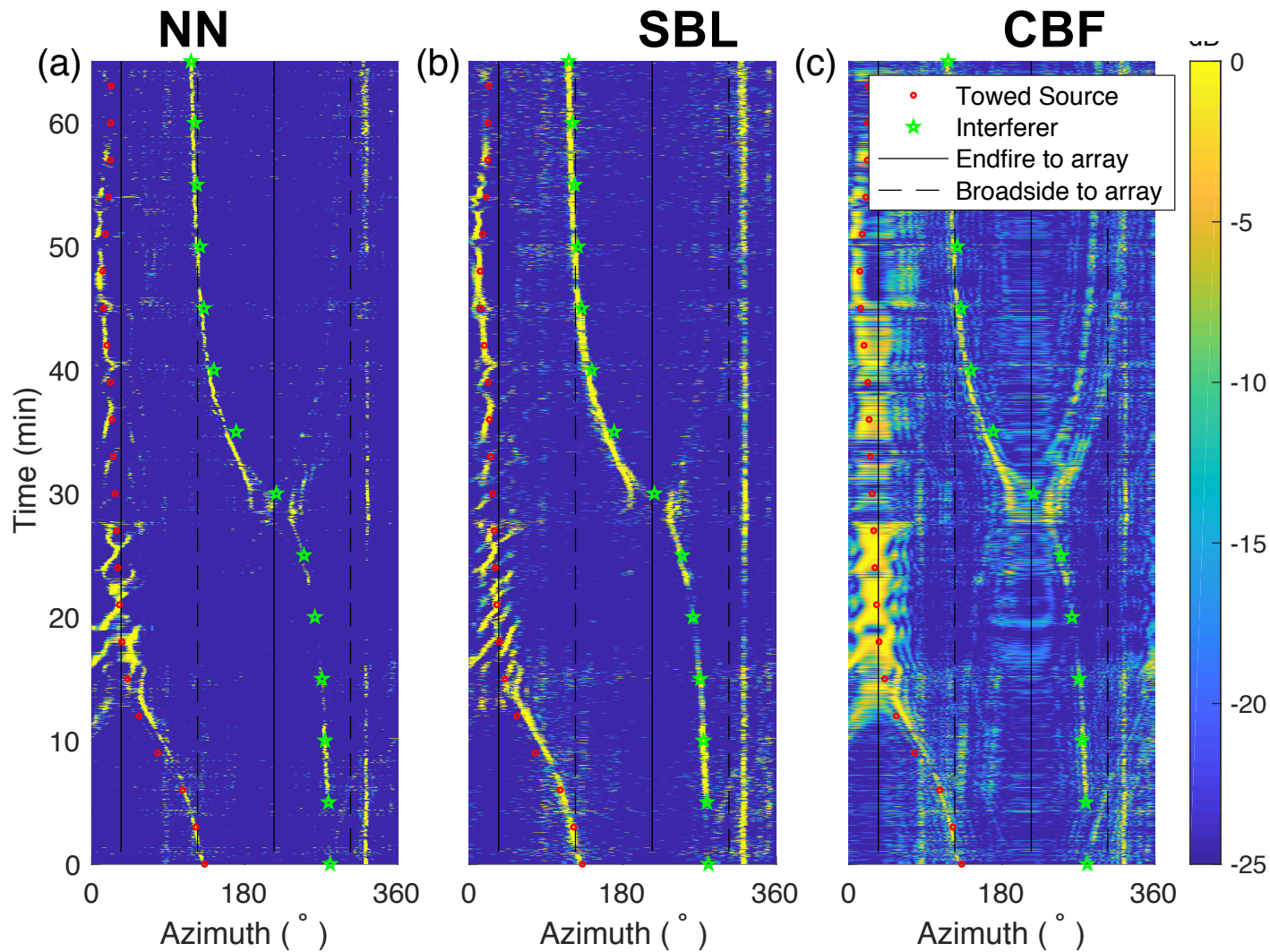
5 layers with 1024 nodes fully connected
20 element array at $\lambda/3$ spacing, searching for 180 DOAs
Coherent sources.



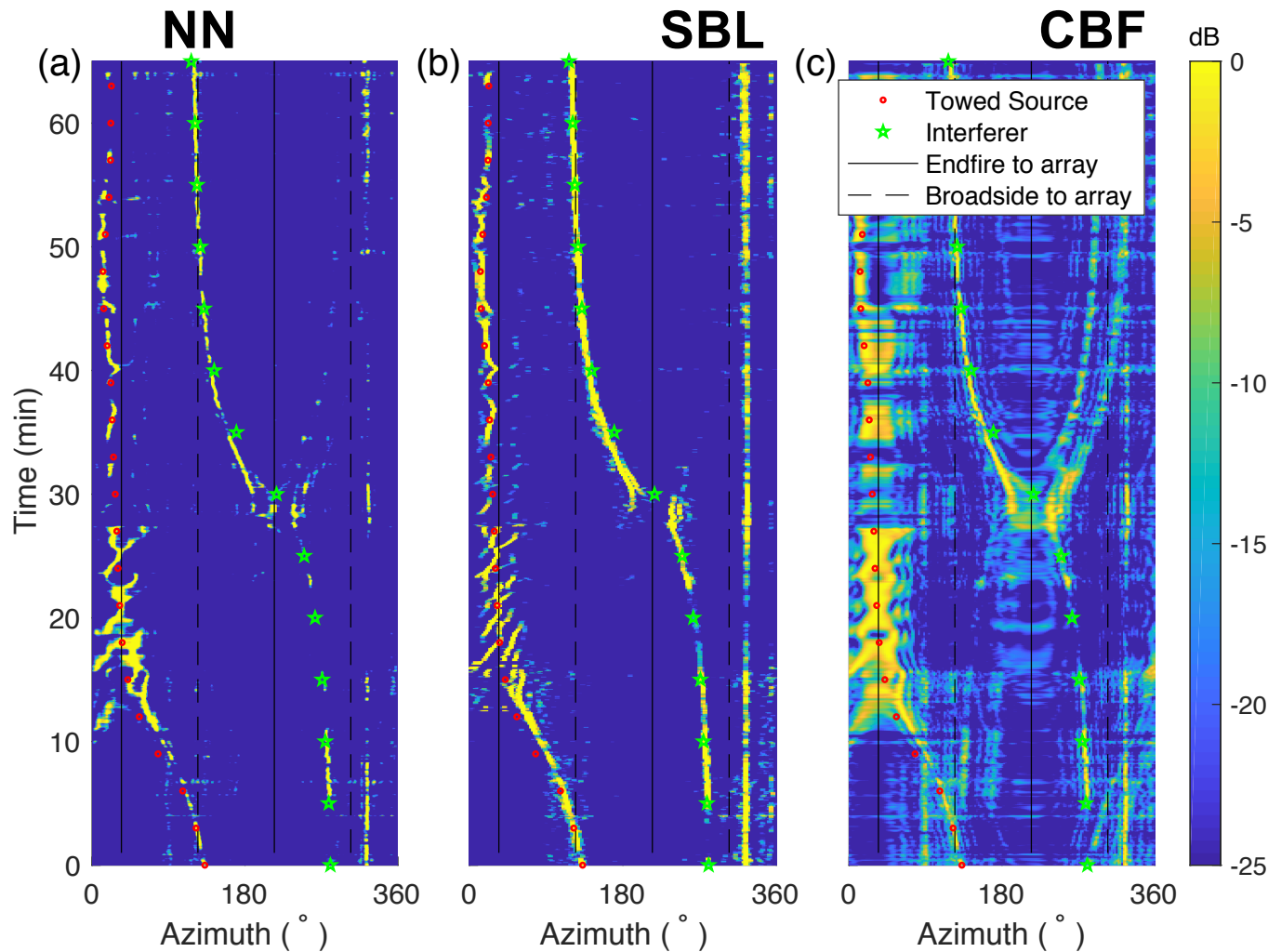Ozanich 2019

# DOA for two sources from SWELLEx96

5 layers with 1024 nodes fully connected
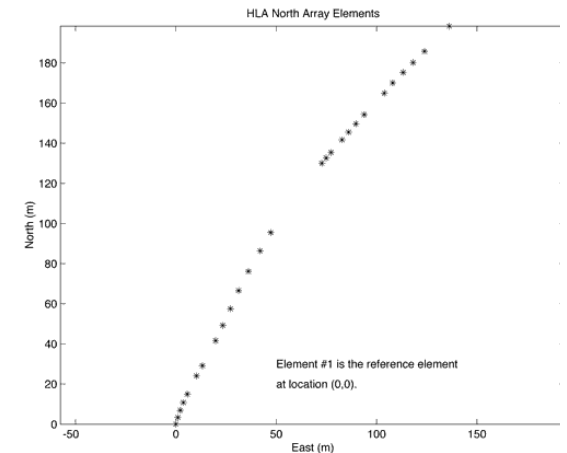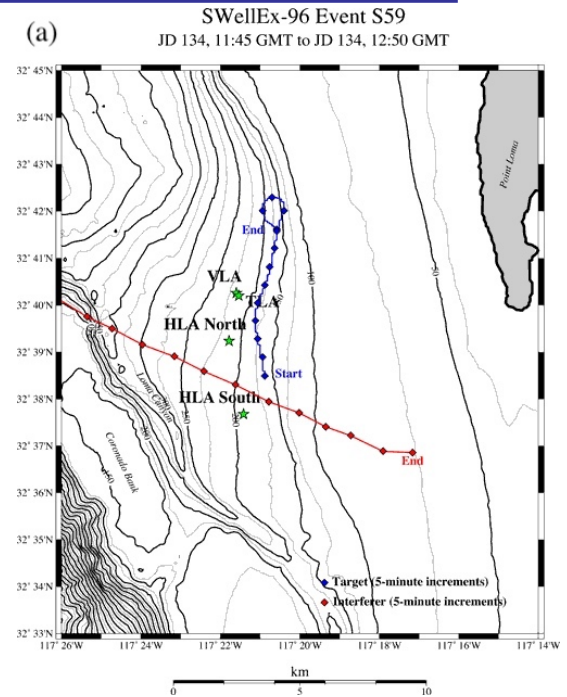One frequency (79 Hz), L=1 snapshot



Ozanich 2019

# DOA for two sources from SW06

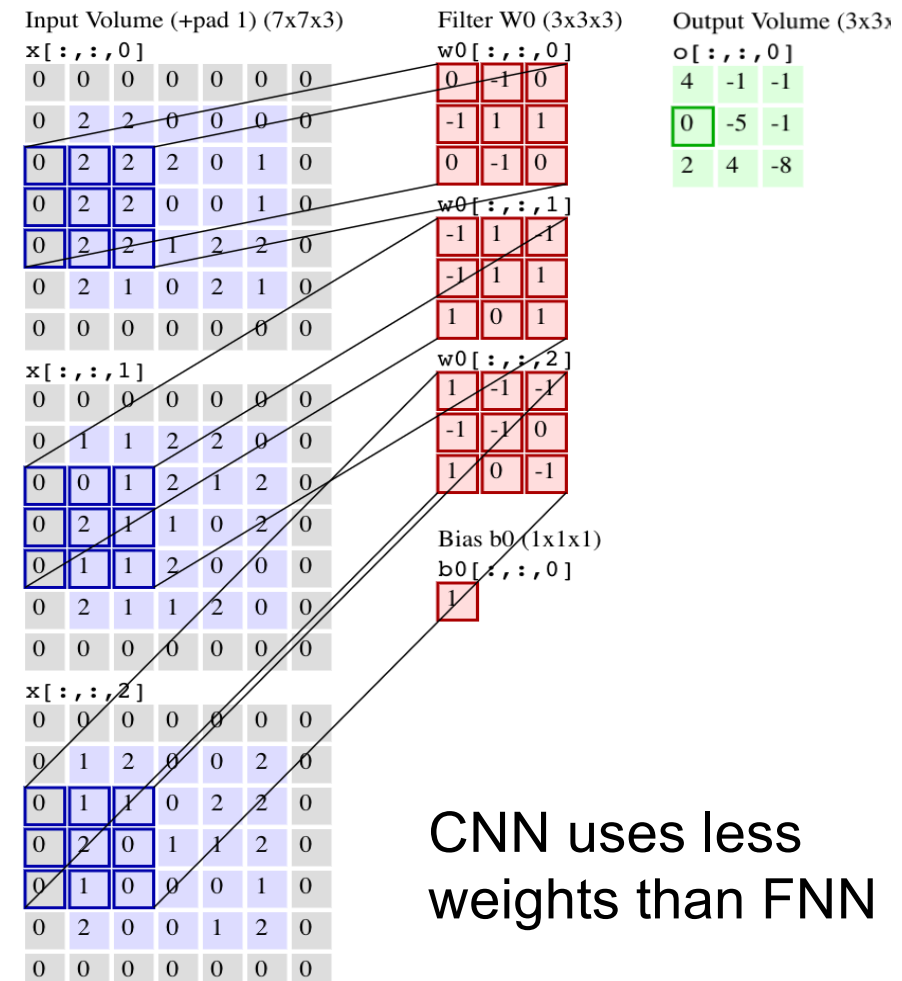5 layers with 1024 nodes fully connected
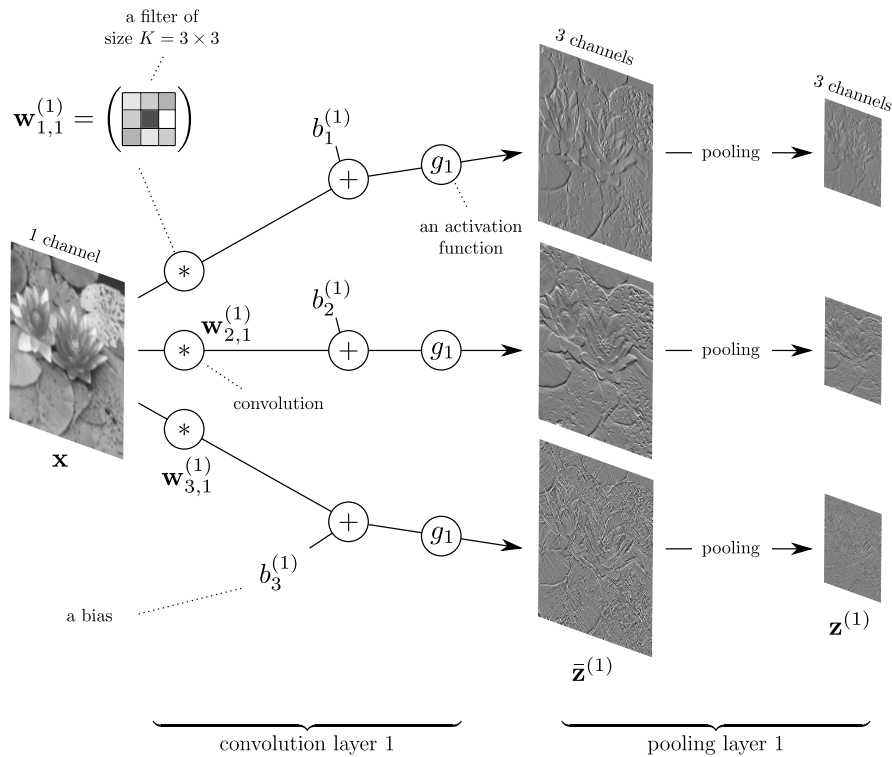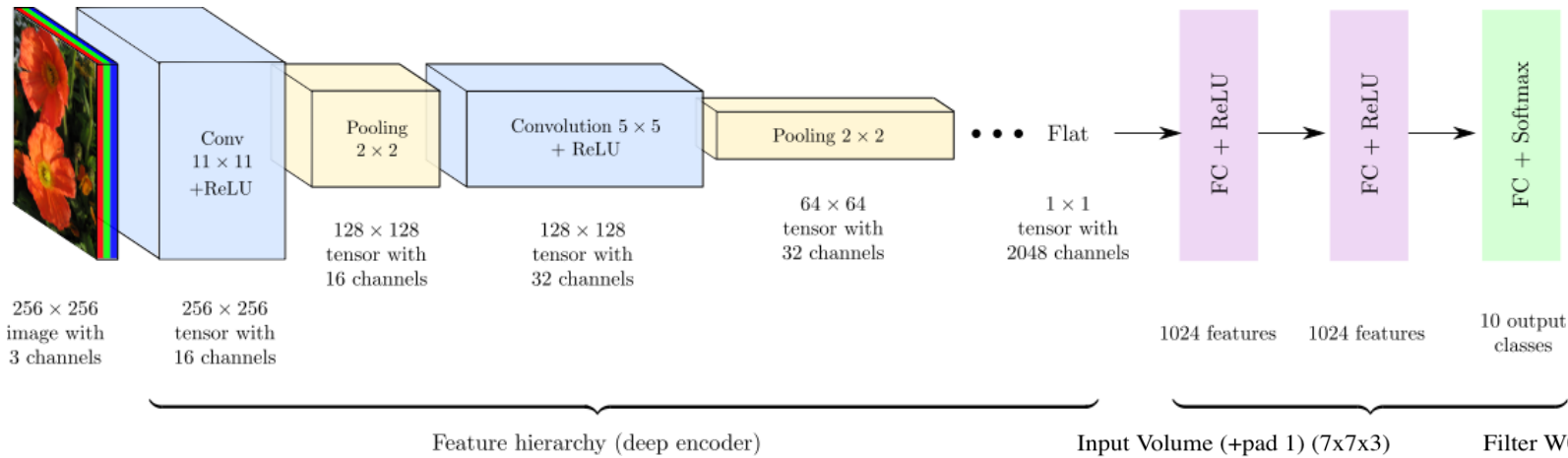One frequency (79 Hz), **L=10** snapshot

**More snapshots give cleaner image**



Ozanich 2019

# Deep Convolutional NN



Conv
$11 \times 11$
+ReLU

Pooling
$2 \times 2$

Convolution $5 \times 5$
+ ReLU

Pooling $2 \times 2$

$\bullet \bullet \bullet$ Flat

FC + ReLU

FC + ReLU

FC + Softmax

$256 \times 256$
image with
3 channels

$256 \times 256$
tensor with
16 channels

$128 \times 128$
tensor with
16 channels

$128 \times 128$
tensor with
32 channels

$64 \times 64$
tensor with
32 channels

$1 \times 1$
tensor with
2048 channels

1024 features

1024 features

10 output
classes

Feature hierarchy (deep encoder)

a filter of
size $K = 3 \times 3$

$\mathbf{w}_{1,1}^{(1)} = $

$b_1^{(1)}$

$g_1$

3 channels

3 channels

pooling

an activation
function

1 channel

$\mathbf{w}_{2,1}^{(1)}$

$b_2^{(1)}$

$g_1$

pooling

convolution

$*$

$\mathbf{w}_{3,1}^{(1)}$

$b_3^{(1)}$

$g_1$

pooling

$\mathbf{x}$

a bias

$\bar{\mathbf{z}}^{(1)}$

$\mathbf{z}^{(1)}$

convolution layer 1

pooling layer 1

Bianco 2019, Niu 2019,

**Input Volume (+pad 1) (7x7x3)**

x[:,:,0]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 2 | 2 | 0 | 0 | 0 | 0 |
| 0 | 2 | 2 | 2 | 0 | 1 | 0 |
| 0 | 2 | 2 | 0 | 0 | 1 | 0 |
| 0 | 2 | 2 | 1 | 2 | 2 | 0 |
| 0 | 2 | 1 | 0 | 2 | 1 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,1]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 2 | 2 | 0 | 0 |
| 0 | 0 | 1 | 2 | 1 | 2 | 0 |
| 0 | 2 | 1 | 1 | 0 | 2 | 0 |
| 0 | 1 | 1 | 2 | 0 | 0 | 0 |
| 0 | 2 | 1 | 1 | 2 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

x[:,:,2]

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 0 | 0 | 2 | 0 |
| 0 | 1 | 1 | 0 | 2 | 2 | 0 |
| 0 | 2 | 0 | 1 | 1 | 2 | 0 |
| 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 2 | 0 | 0 | 1 | 2 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

**Filter W0 (3x3x3)**

w0[:,:,0]

| 0 | -1 | 0 |
|---|---|---|
| -1 | 1 | 1 |
| 0 | -1 | 0 |

w0[:,:,1]

| -1 | 1 | -1 |
|---|---|---|
| -1 | 1 | 1 |
| 1 | 0 | 1 |

w0[:,:,2]

| 1 | -1 | -1 |
|---|---|---|
| -1 | -1 | 0 |
| 1 | 0 | -1 |

Bias b0 (1x1x1)
b0[:,:,0]

| 1 |
|---|

**Output Volume (3x3x)**

o[:,:,0]

| 4 | -1 | -1 |
|---|---|---|
| 0 | -5 | -1 |
| 2 | 4 | -8 |

## CNN uses less weights than FNN
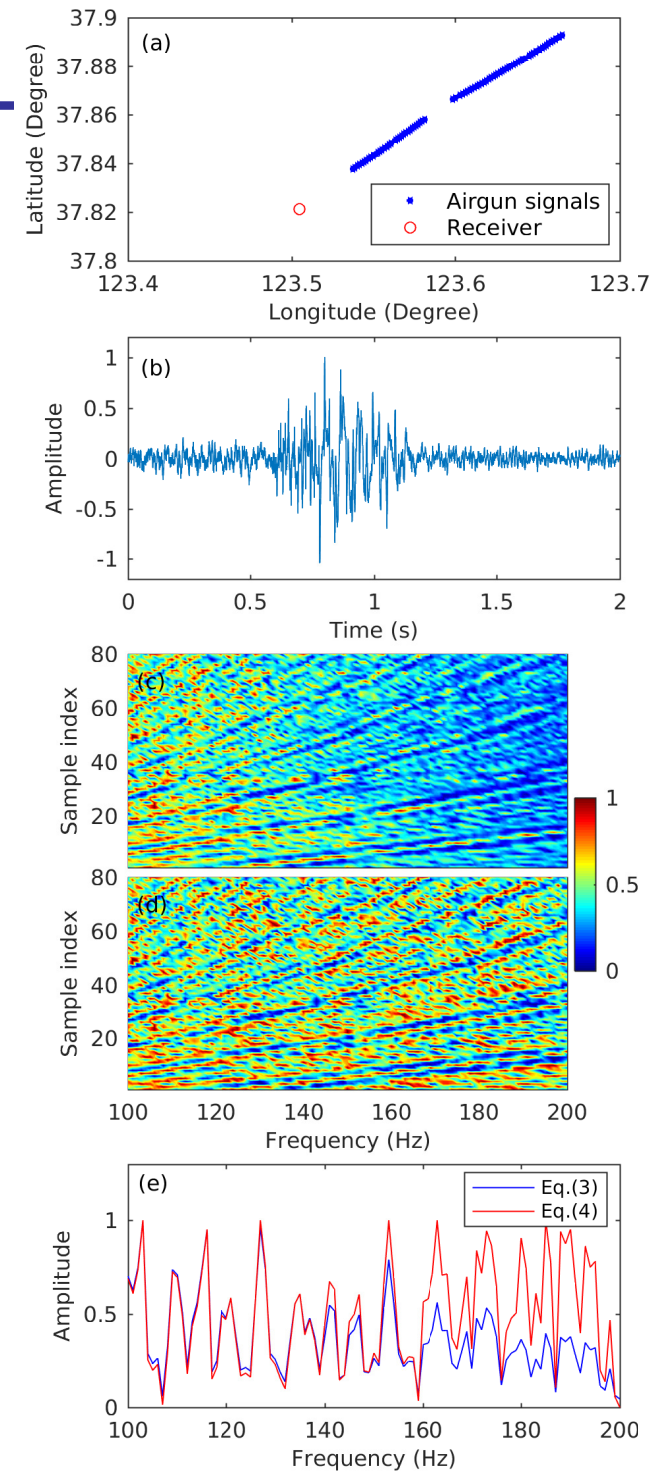
# Magnitude only localization



Single receiver,
3-16 km from source
Multi-frequency 100-200 Hz,
magnitude only

Much less input as sample covariance matrix is not needed. Magnitude is averaged directly
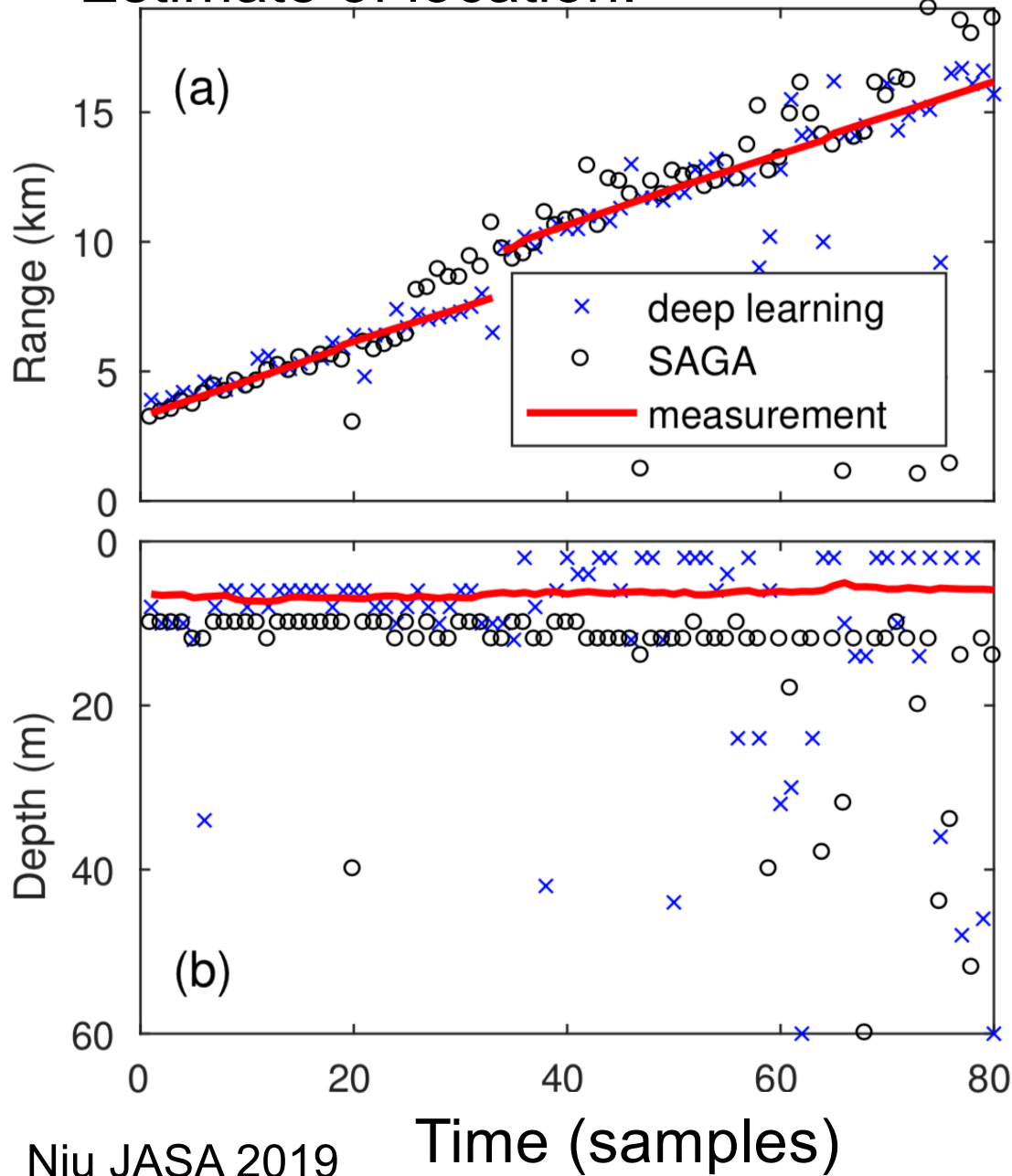
SAGA, multi frequency objective function

$$\phi_F(\Theta) = 1 - \frac{|\sum_{f=1}^{F} \hat{\mathbf{p}}(f)\hat{\mathbf{q}}(f,\Theta)|^2}{\sum_{f=1}^{F} |\hat{\mathbf{p}}(f)|^2 \sum_{f=1}^{F} |\hat{\mathbf{q}}(f,\Theta)|^2},$$

$\widehat{p}$ and $\widehat{q}$ are magnitudes

Niu 2019

# ML and SAGA ranging

## Estimate of location:



## Statistics of location



## Does ML **beat** SAGA?

Niu JASA 2019

Time (samples)

# Graph Signal Processing for locating a source

**Location 2:** Otis Redding - "Hard to handle"

f = 750 Hz

30-microphone array

*i*    *j*

Spectral coherence

$$\hat{C}_{ij}(f) = \frac{1}{N}\sum_{t=1}^{N} X_i(f,t) \cdot \bar{X}_j(f,t)$$

*(Normalization: |X(f,t)|²=1)*

**Location 1:** Prince - "Sign o' the times"

Riahi 2017

25

Statistically significant entries
=> **Connectivity matrix**



Rec. no.

Graph with 30 nodes

- Each sensor is a **node** in the graph.
- If **nodes** $i$ and $j$ are significantly correlated $|C_{ij}|>\xi$, then they share an **edge**.
- A **subgraph** has high spatial coherence across a subarray (=> likely a source nearby).

**Connected subgraphs:**

**5 nodes and 9 edges**

**8 nodes and 20 edges**

Riahi 2017

# Graph clustering for localization within a sensor array

Peter Gerstoft and Nima Riahi,   **noiselab.ucsd.edu**
Christoph Mecklenbrauker, TU Wien

Based on paper: Riahi and Gerstoft, Signal Processing, 2017

March 5—12, 2011:  3TB, 5200 Stations in Long Beach, California

05:53:09--05:55:06h

47 Hz

UTME [km]

**Helicopter rotor noise (seismo-acoustic coupling)**

Several peaks consistent with helicopter rotor harmonics (20-100 Hz).

Doppler shift
$f_{high}/f_{low}=(v_0+v)/(v_0-v)\approx1.4$ i.e. $v\approx250$ km/h

Speed over ground 7km/2min=210km/h

✓ Rotor frequencies

✓ Doppler frequency shift

✓ Movement in map

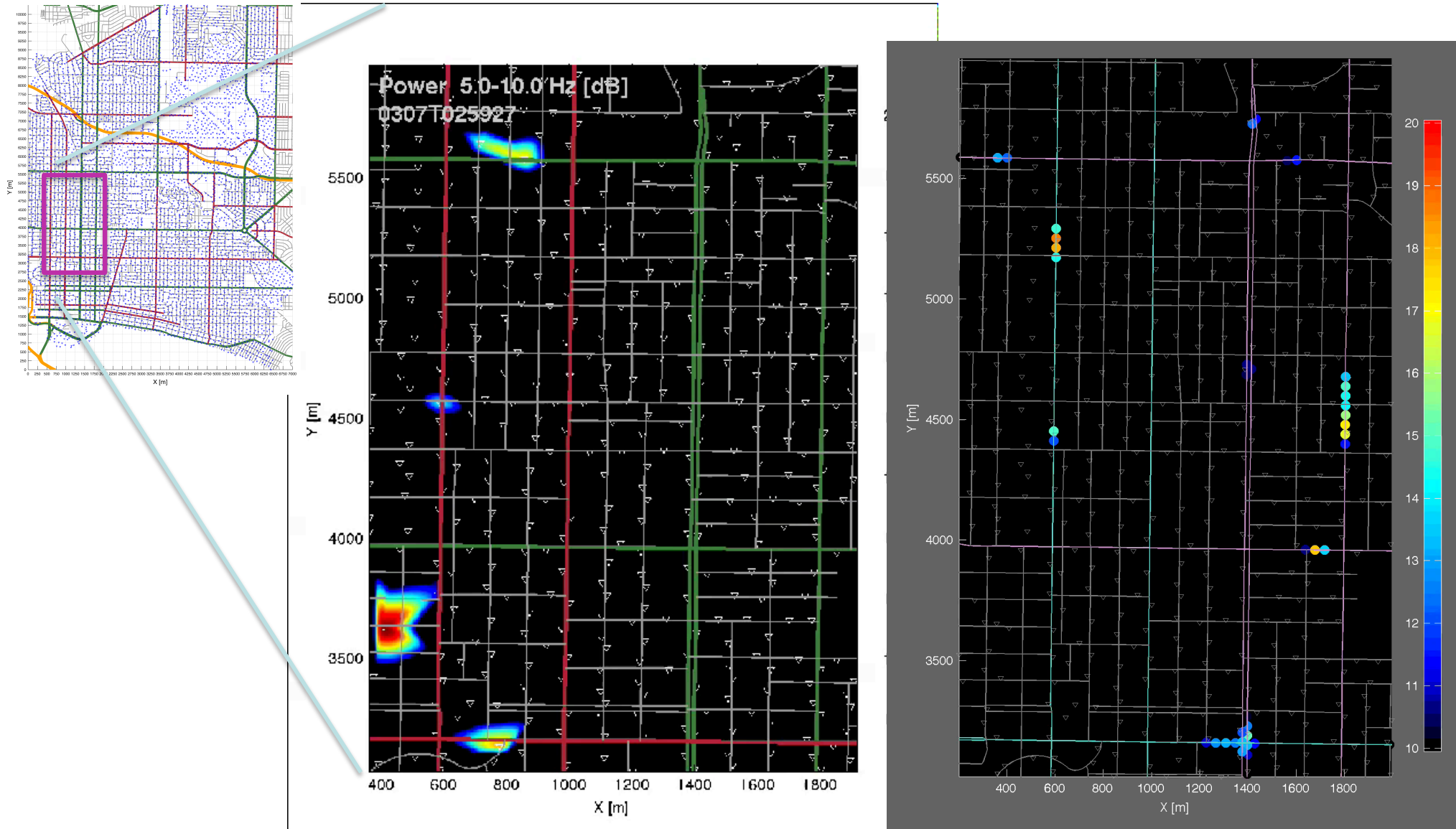Aliased energy

12 Hz

Freq [Hz]

time [MM:SS]

dB

28

# Clusters on March 10



**10-19Hz**

airport

**pump jacks and drill rigs
2: Pumping facility**
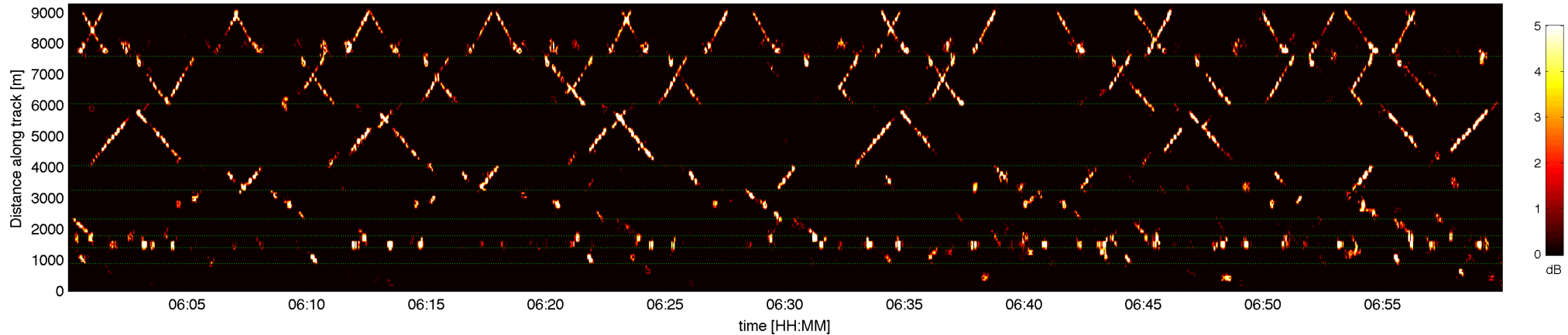
**Long Beach light rail
(Blue Line Metro)**

Based on 9400 time windows x 10 frequency bins.

Each dot is the center of a cluster. 90% of the clusters cover <1.5% of the area.

Few false detections

Riahi, Gerstoft, Signal Processing 2017

## 5200 element Long Beach array (Dan Hollis)



Riahi, Gerstoft, The seismic traffic footprint: Tracking trains, aircraft, and cars seismically, GRL 2015

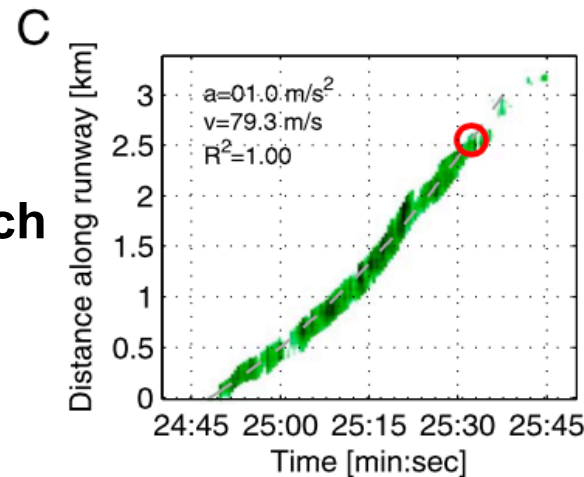# Noise Tracking of Cars/Trains/Airplanes
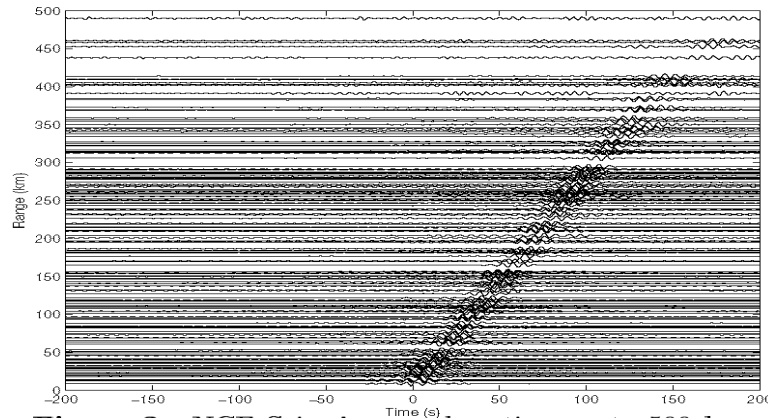
**March 7th, 6-7am, rush hour, Blue Line**

**Accelerating airplane on Long Beach airport runway, moving northwest and taking off at about 120 mi/h.**

Riahi, Gerstoft, GRL 2015

# Travel time tomography

### Travel times from noise cross-correlations



distance = speed x time

slowness = 1/speed

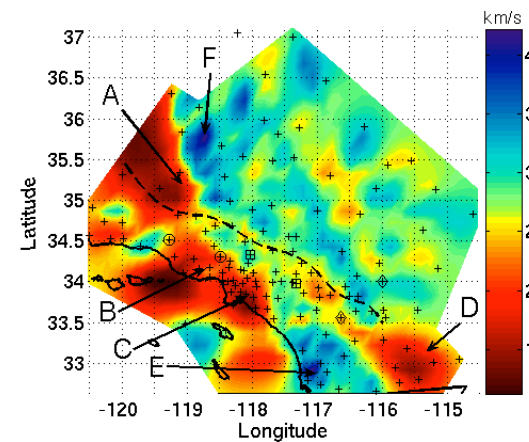- Task: Given travel times, estimate regional phase speed distribution

$$d = Am + n,$$

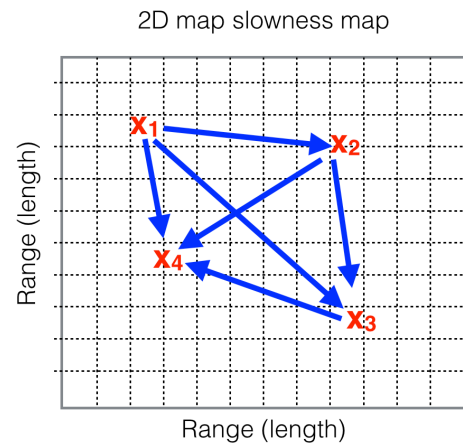$d$: M travel times

$A$: "Tomography matrix": ray paths through the discretized map

$m$: N-pixel slowness image

Slowness map and measurements
- stations in red
- rays in blue
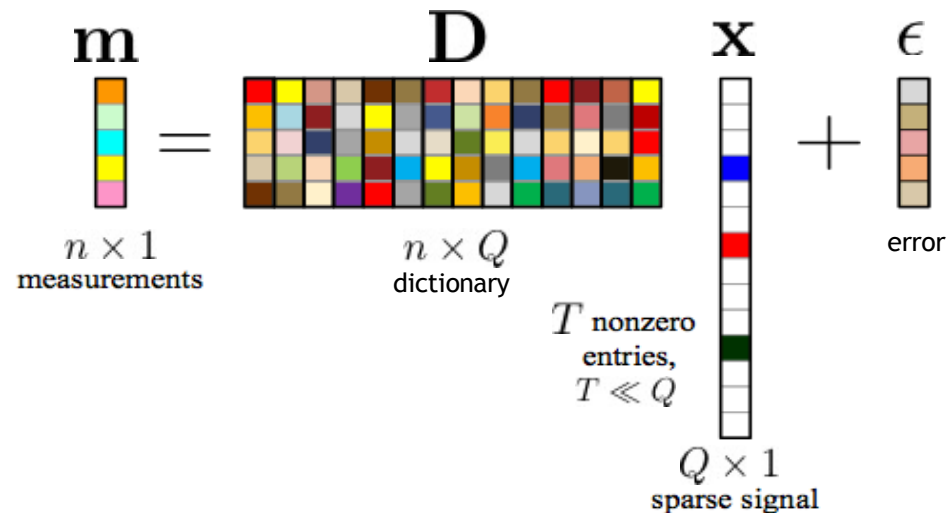
### 2D map slowness map



Low Velocity Region~Sedimentary basins A: San Joaquin, B: Ventura, C: L.A., D: Salton Sea, E: Peninsular range, F: Sierra Nevada

# Sparse models and dictionaries

- Sparse modeling assumes each signal model can be reconstructed from a few vectors from a large set of vectors, called a dictionary **D**
- Adds auxiliary sparse model to measurement model

$$\mathbf{d} = \mathbf{A}\mathbf{m} + \mathbf{n}, \ \mathbf{m} \approx \mathbf{D}\mathbf{x} \text{ and } |\mathbf{x}| \ll Q$$

- Optimization changes from estimating **m** to estimating sparse coefficients **x**



$\mathbf{m}$    $\mathbf{D}$    $\mathbf{x}$    $\epsilon$

$n \times 1$ measurements    $n \times Q$ dictionary    error

$T$ nonzero entries, $T \ll Q$
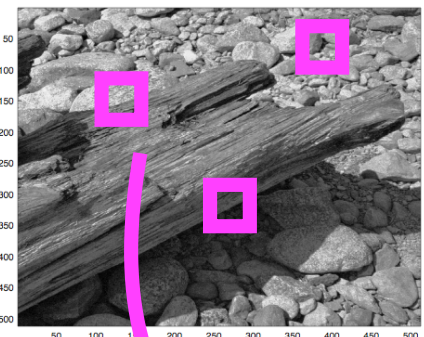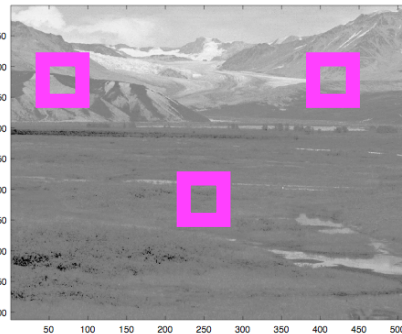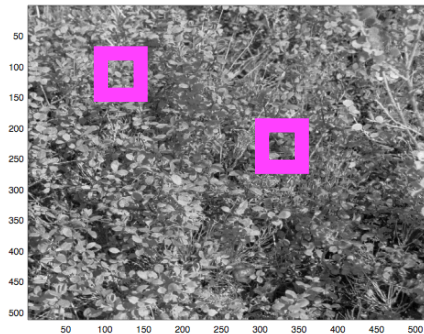
$Q \times 1$ sparse signal

- Sparse objective: $\min_{\mathbf{x}} \|\mathbf{A}\mathbf{D}\mathbf{x} - \mathbf{d}\|_2$ subject to $\|\mathbf{x}\|_0 \leq T$
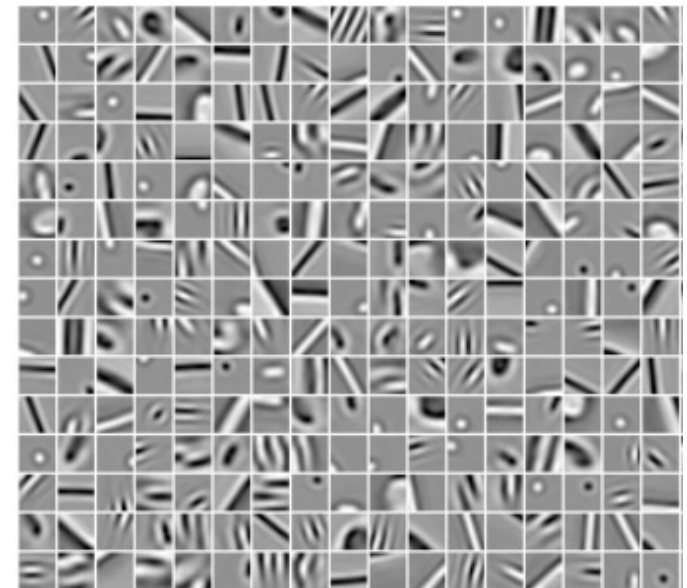
# Dictionary learning and sparsity

- Dictionary learning obtains "optimal" sparse modeling dictionaries directly from data

- Dictionary learning was developed in neuroscience (a.k.a. sparse coding) to help understand mammalian visual cortex structure

- Assumes (1) <u>Redundancy in data:</u> image patches are repetitions of a few elemental shapes; and (2) <u>Sparsity:</u> each patch is represented with few atoms from dictionary

"Natural" images, patches shown in **magenta**

Learn dictionary $\mathbf{D}$ describing $\mathbf{Y} = [\mathbf{y}_1, ..., \mathbf{y}_I]$

- Each patch is signal $\mathbf{y}_i$
- Set of all patches $\mathbf{Y} = [\mathbf{y}_1, ..., \mathbf{y}_I]$
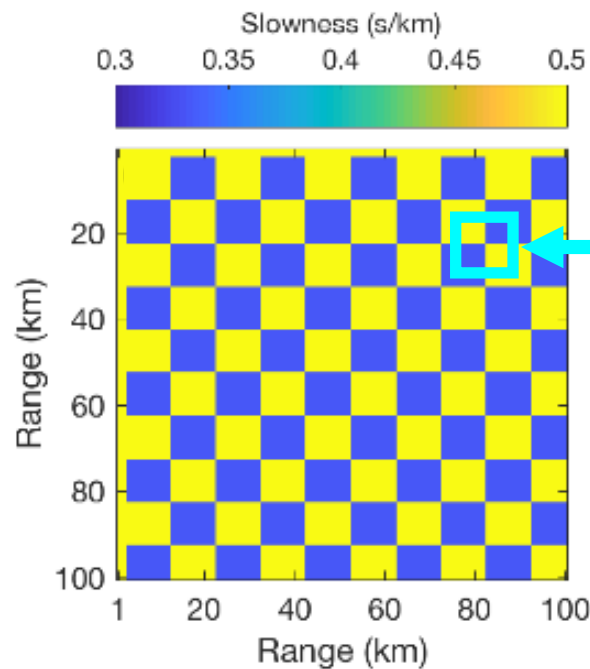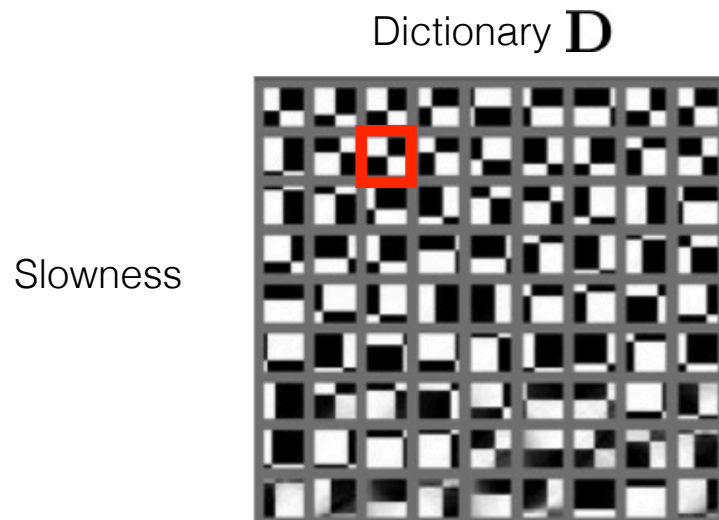
Olshausen 2009

Bianco 2018, 2019

Sparse model for patch $\mathbf{y}_i$ composed of few atoms from $\mathbf{D}$

$$\widehat{\mathbf{x}}_i = \arg\min_{\mathbf{x}_i} \|\mathbf{y}_i - \mathbf{D}\mathbf{x}_i\|_2 \quad \text{subject to} \quad \|\mathbf{x}_i\|_0 \leq T$$

$$\mathbf{y} = \boxed{} = \boxed{} x_1 + \boxed{} x_2 + ...$$

# Checkerboard dictionary example

Dictionary $\mathbf{D}$

Slowness



Slowness (s/km)

$$\mathbf{y} = \mathbf{R}_i \mathbf{s} = \mathbf{D}\mathbf{x}_i$$
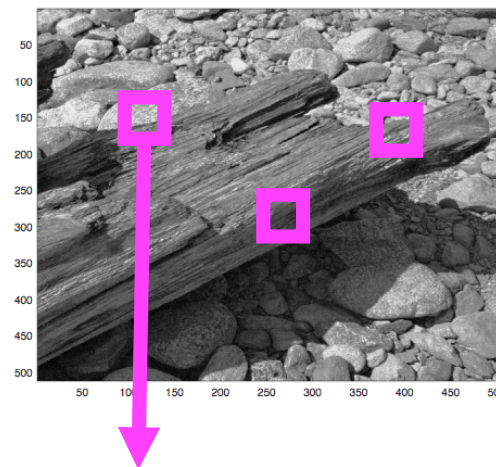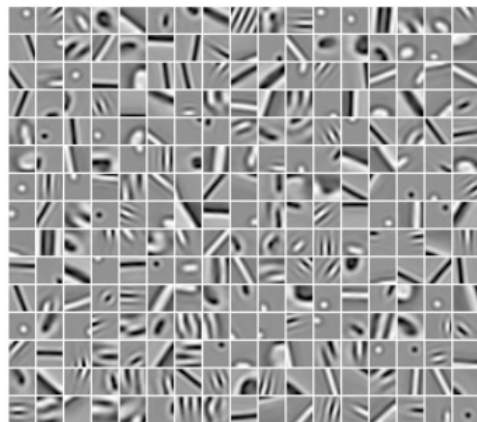
$$\mathbf{R}_i \mathbf{s} = \ \blacksquare \ x_i$$

10x10 pixel patches

$$\mathbf{D} \in \mathbb{R}^{n \times Q}$$

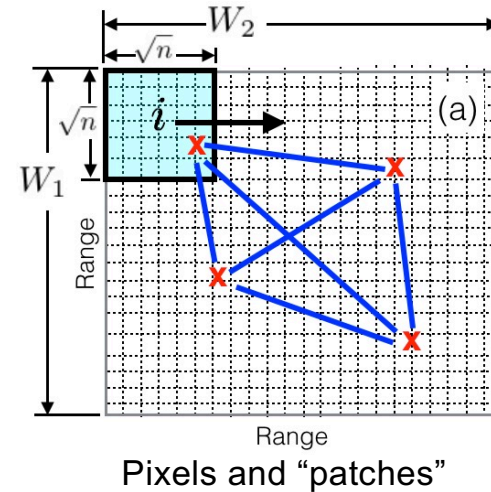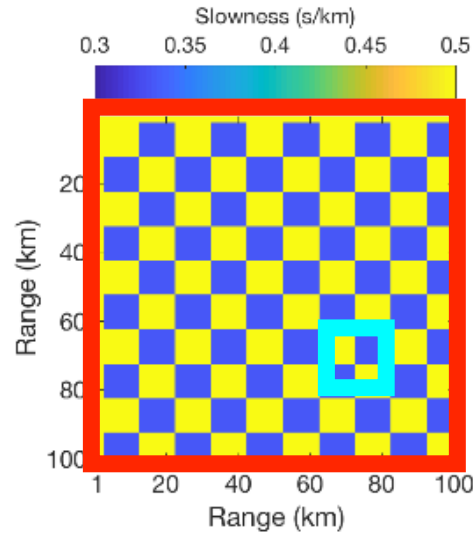$$\mathbf{R}_i \in \{0, 1\}^{n \times N}$$

Natural image

$$\widehat{\mathbf{x}}_i = \arg\min_{\mathbf{x}_i} \|\mathbf{y}_i - \mathbf{D}\mathbf{x}_i\|_2 \quad \text{subject to} \quad \|\mathbf{x}_i\|_0 \leq T$$

Bianco 2018, 2019

$$\mathbf{y} = \ \blacksquare \ = \ \blacksquare \ x_1 + \ \blacksquare \ x_2 + \ldots$$

# LST slowness image and sampling



Slowness map and measurements
- **stations in red**
- **rays in blue**

Pixels and "patches"

Slowness map and sampling:

- Discrete slowness map $N=W_1 \times W_2$ pixels
- $I$ overlapping $\sqrt{n} \times \sqrt{n}$ pixel patches
- $M$ straight-ray paths

Tomography matrix (straight ray)

$$\mathbf{A} \in \mathbb{R}^{M \times N}$$

Slowness dictionary

$$\mathbf{D} \in \mathbb{R}^{n \times Q}$$
$$Q \ll I$$

"Local" model

$$\widehat{\mathbf{x}}_i = \arg\min_{\mathbf{x}_i} \|\mathbf{R}_i \mathbf{s}_s - \mathbf{D}\mathbf{x}_i\|_2^2 \text{ subject to } \|\mathbf{x}_i\|_0 = T$$

"Global" model

$$\mathbf{t} = \mathbf{A}\mathbf{s}_g + \epsilon \qquad \widehat{\mathbf{s}}_g = \arg\min_{\mathbf{s}_g} \|\mathbf{t} - \mathbf{A}\mathbf{s}_g\|_2^2 + \lambda_1 \|\mathbf{s}_g - \mathbf{s}_s\|_2^2,$$
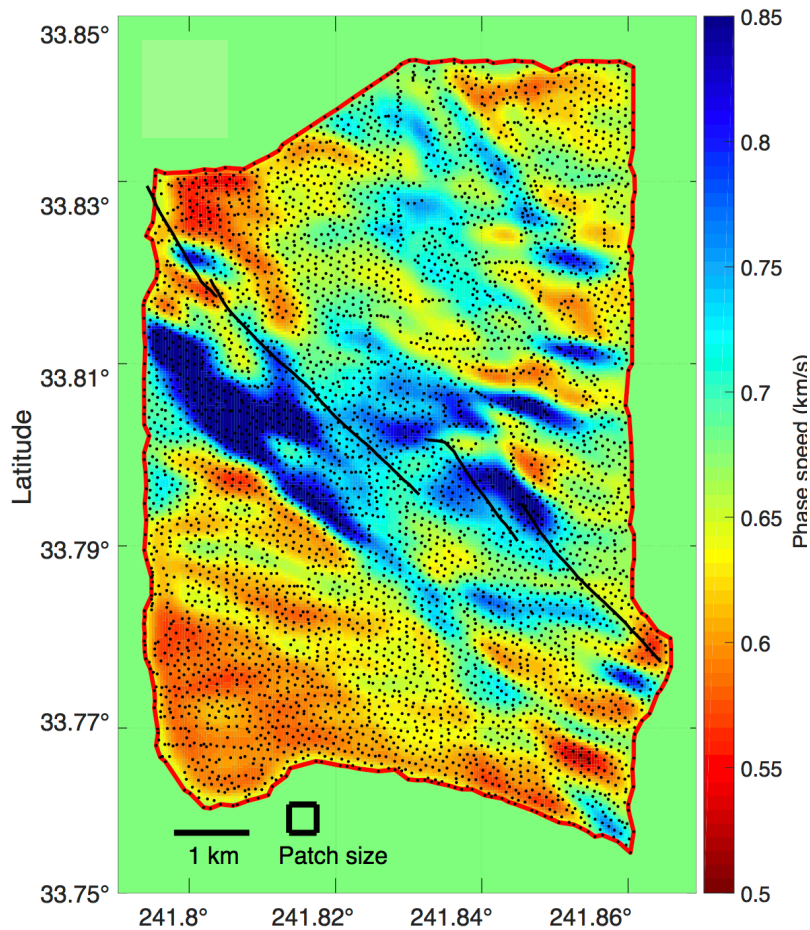
Bayesian formulation

# LST versus conventional tomography

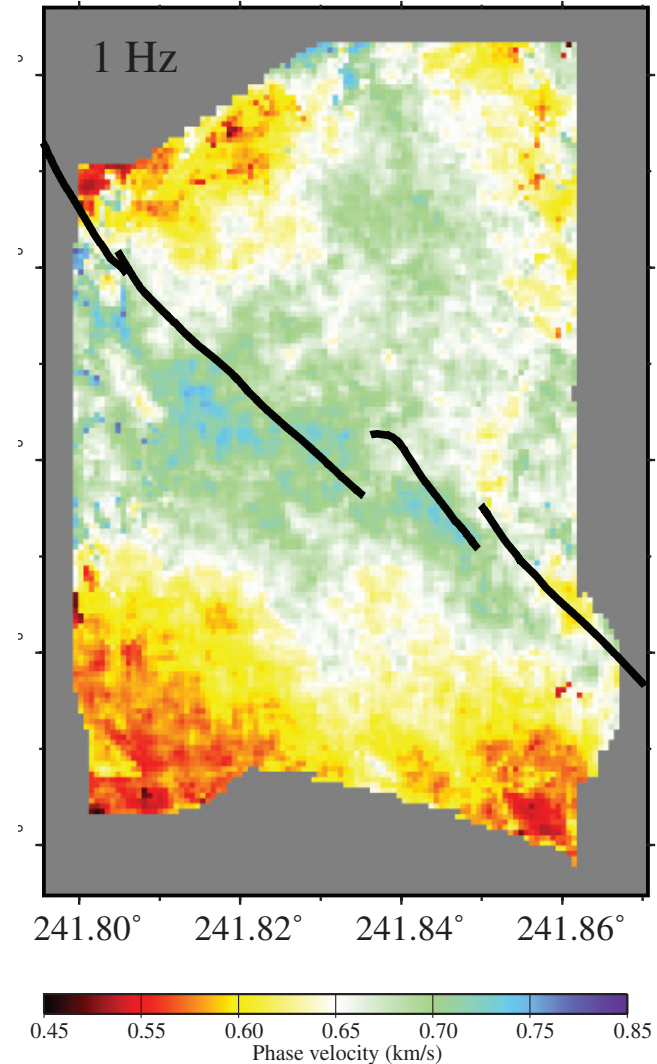Both use same travel times (from Fan-Chi Lin),   <span style="color:red">unsupervised</span>

**LST 3 mill rays**

**Fan-Chi Lin, Geophysics, 8mill Rays**



W$_1$=200, W$_2$=300 pixels
n=100, Q=200, T=1
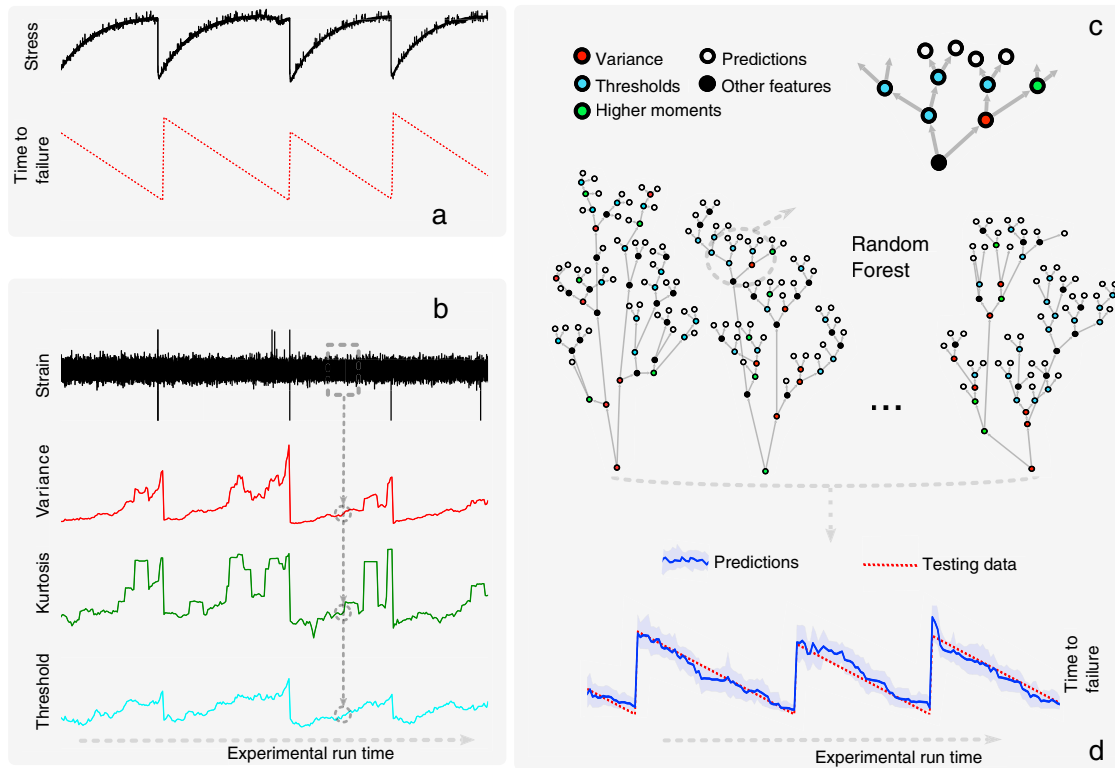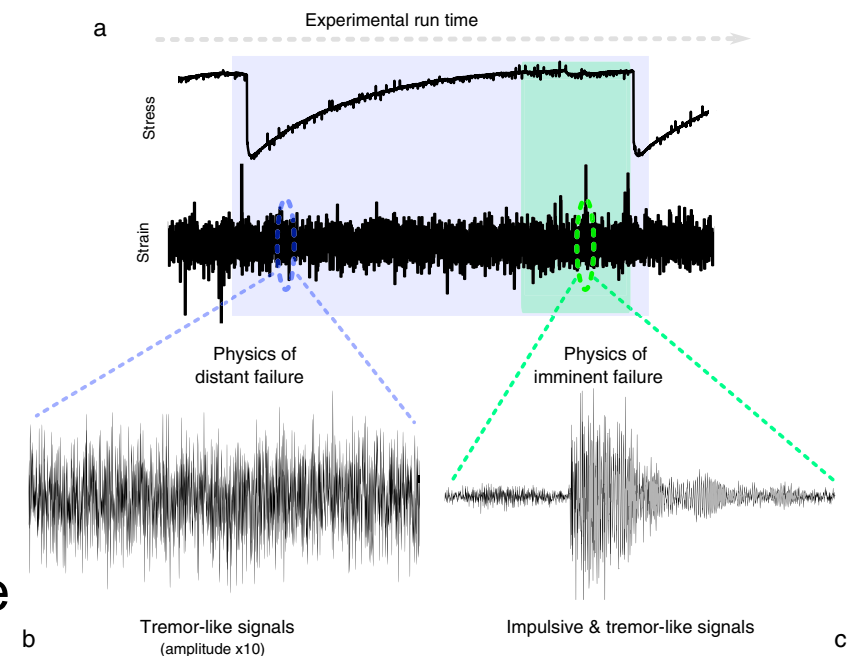
Bianco 2018, 2019

# Predicting Earthquakes in Laboratory

- Kaggle competition.



Once they found a ML that could predict lab-EQ, they also could see the feature.



ML gives little or no insight into the model. We want the ML algorithm to provide a line of reasoning together with the calculated result. Not just the outcome of Bayes formalism.

**=> That will come**

## First principles   vs   Data driven

| | First principles | Data driven |
|---|---|---|
| Data | Small data | Big data to train |
| Domain expertise | High reliance on domain expertise | Results with little domain knowledge |
| Fidelity/ Robustness | Universal link can handle non-linear complex relations | Limited by the range of values spanned by training data |
| Adaptability | Complex and time consuming derivation to use new relations | Rapidly adapt to new problems |
| Interpretability | Parameters are physical! | Physically agnostic, limited by the rigidity of the functional form |
| Perceived Importance. | **Geophys**    **SignalProc** | **Peter**    **Googl** |

# Summary

- Machine learning, big data, data science, artificial intelligence are similar.

- **Data science** has lots of opportunities in **geophysics**.

- Neural networks is one method. Similar methods are Support Vector Machines (SVM) and Random Forrest (RF). Use the latter for a first test.

- **Unsupervised learning** is more challenging than supervised learning

- We need explainable artificial intelligence. We want the ML algorithm to provide a line of reasoning together with the calculated result / fit / decision.

**Actions: Download  ML JASA review**

- TRY http://playground.tensorflow.org

Can ML

- Replace CTBTO processing chain?

- Discover PDE (Partial differential equation) in video?

- Find sea mines?

- Design metamaterials?

- Predict earthquakes?

- Replace 50 years of array processing

- Source location in the ocean waveguide w/o training.

# FINITO