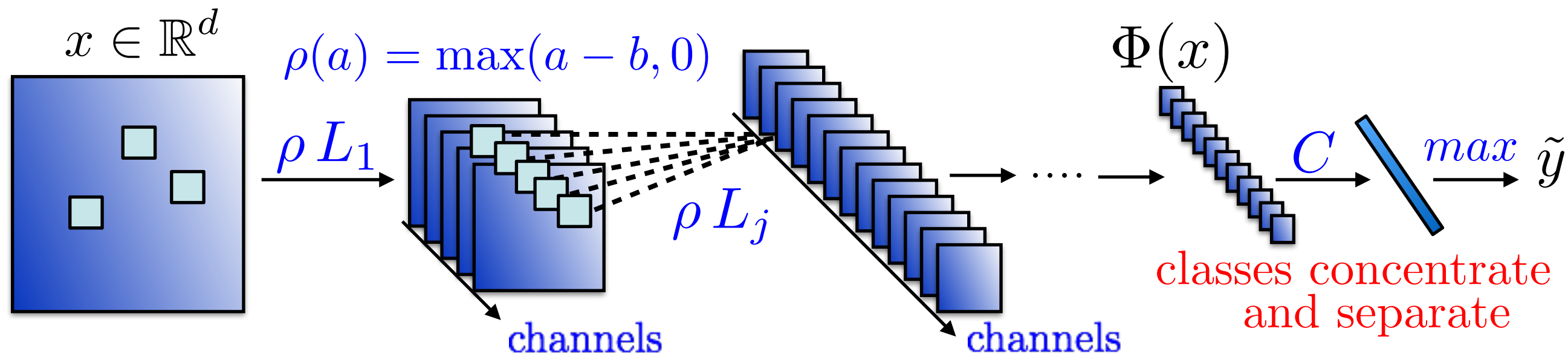# A Harmonic Analysis View of

# Deep Network Theory

*Stéphane Mallat*

**Flatiron Institute, CCM**

**Collège de France, ENS Paris**

# A View of Deep Network Theory

Classification with deep convolutional networks:



$x \in \mathbb{R}^d$

$\rho(a) = \max(a - b, 0)$

$\rho\, L_1$

$\rho\, L_j$

channels

$\Phi(x)$

$C$

$max$ $\tilde{y}$

classes concentrate and separate

channels

- Surprisingly good generalisation properties: not understood
- Issues of robustness and validation in applications: *transport, medecine, sciences...*
- Opportunity for new maths and science theories

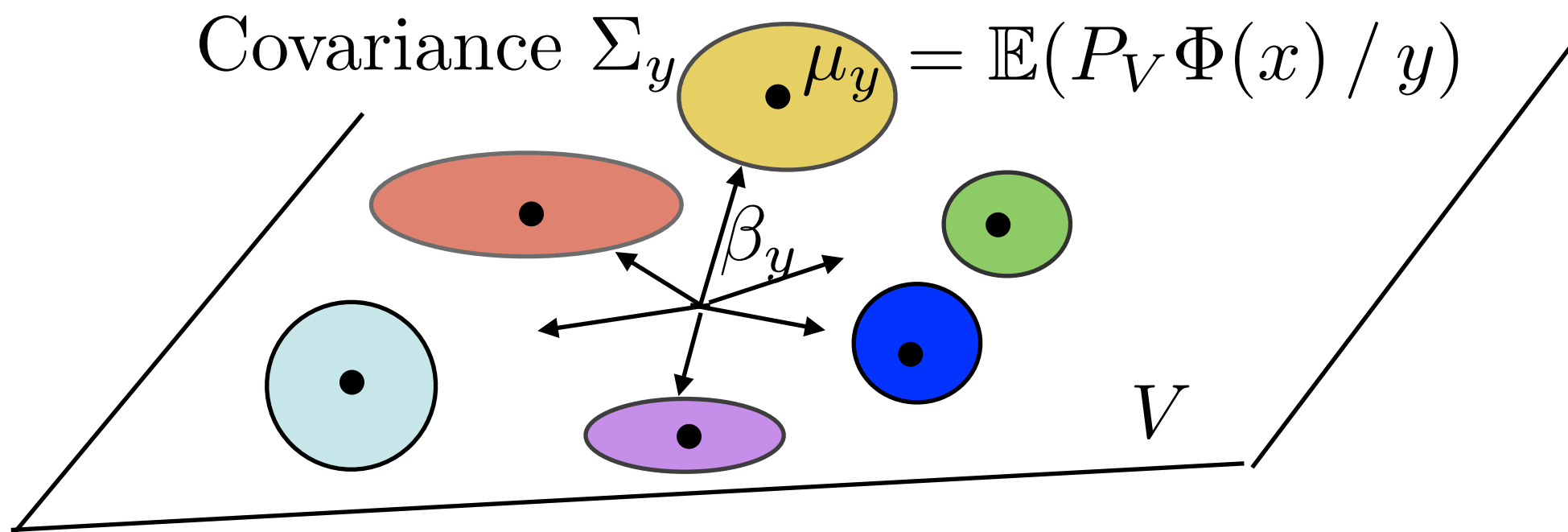*Weekly working seminar with university collaborations (Joan Bruna)*

Tuesdays 11am-12am EST (CCM Web page)

November 10th: *David Donoho* Neural Collapse

Linear classifier: $\tilde{y} = \arg_y \max \langle \Phi(x), \beta_y \rangle + \alpha_y$

Only depends on the projection of $\Phi(x)$ on $V = Vect\{\beta_y\}_y$:

Covariance $\Sigma_y$ $\mu_y = \mathbb{E}(P_V \Phi(x) / y)$



$\beta_y$

$V$

- $P_V \Phi(x)$ must have separated class means $\mu_y$:

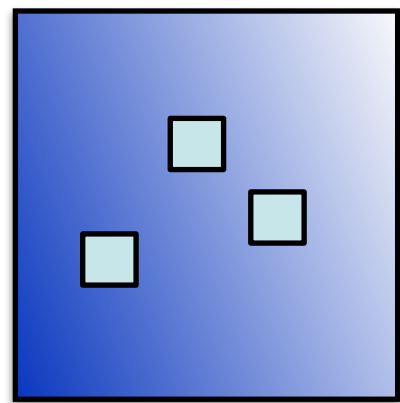Fisher Ratio: $\text{Trace}(\Sigma_W^{-1} \Sigma_B) \xrightarrow[\text{training}]{\text{Neural collapse}} \infty$

V. Papyan
X.Y. Han
D. Donoho

with $\Sigma_B = Ave_y (\mu_y - \overline{\mu})(\mu_y - \overline{\mu})^T$ and $\Sigma_W = Ave_y \Sigma_y$.

What $\Phi(x)$ achieves this concentration/separation ?
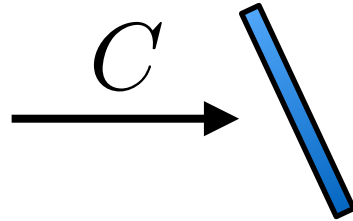
# Tight Frame Contraction

*John Zarka, Florentin Guth*



$$x \in \mathbb{R}^d \xrightarrow{\rho F} \Phi = \rho F \xrightarrow{C} C\Phi(x) = \left( \langle \Phi(x), \beta_y \rangle \right)_y$$

$\Phi = \rho F$ increase dimension

$S = C\,\rho\,F$ : 2 layer network with no bias

Tight frame: $F^T F = Id$,

contraction: $|\rho(a) - \rho(a')| \le |a - a'|$ "Stein shrinking estimation"
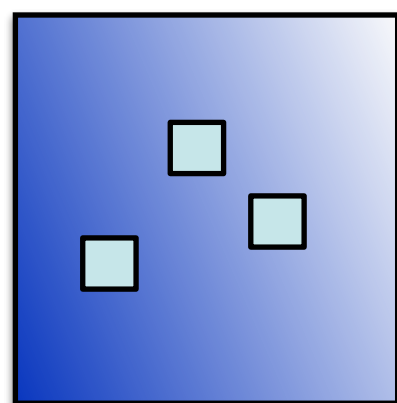
$$\Rightarrow \quad \|\Phi(x) - \Phi(x')\| \le \|x - x'\| \; : \; \text{contraction}$$

Contractions with a fixed global bias $b_0$:

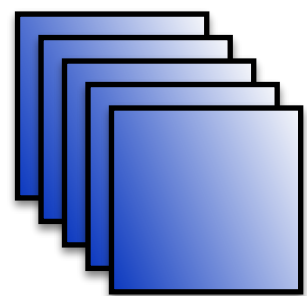Soft-Thresh. $\rho(a) = \text{sign}(a) \max(|a| - b_0, 0)$ *shrinks amplitude* for noise removal

ReLu $\rho(a) = \max(a - b_0, 0)$ *shrinks amplitude and sign*
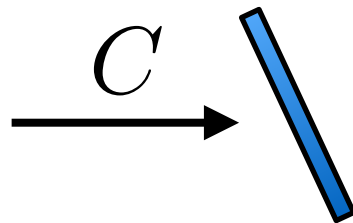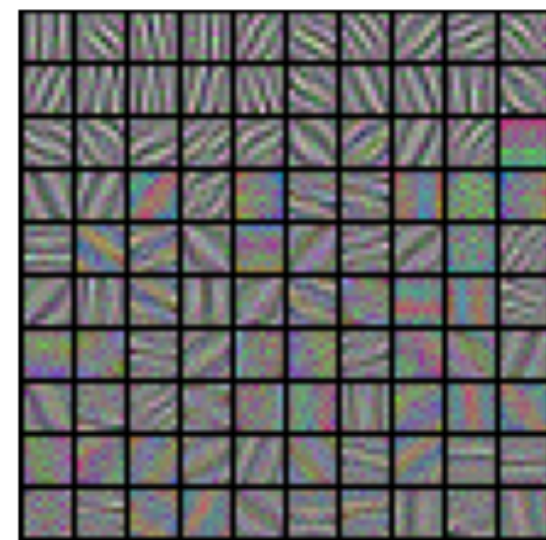
# Tight Frame Contraction



$x \in \mathbb{R}^d$

$\rho F$
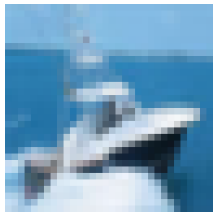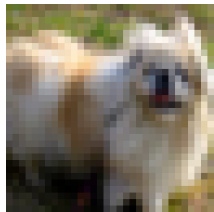
$\Phi = \rho F$
increase dimension

$C$

$S$

Filters of $F$ for CIFAR

- SGD optimisation

| $\Phi(x)$ | | $x$ | Soft $\rho F x$ | ReLu $\rho F x$ |
|---|---|---|---|---|
| MNIST | Error | 7.4% | 1.4% | 1.4% |
| | Fisher | 20 | 60 | 60 |
| CIFAR | Error | 60% | 39% | 28% |
| | Fisher | 7 | 12 | 15 |

- A soft-thresholding $\rho$ can reduce within class variance and preserve class means $\mu_y$ if $Fx$ is sufficiently sparse. *(Donoho Johnstone)* A ReLu $\rho$ also modifies class means.

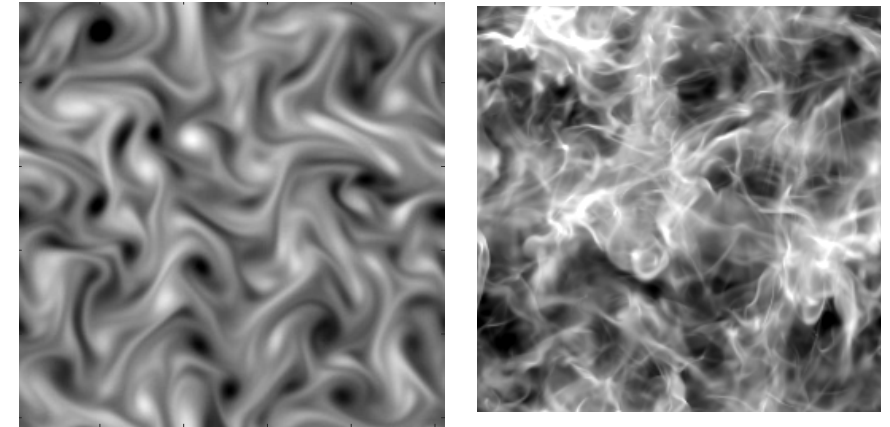Do we need to learn the tight frame $F$ ?

# Overview



**I- Concentration in Statistical Physics:**

- Models of non-Gaussian processes

  <span style="color:blue">Turbulences:</span>

- Wavelet separation and ReLU: scales, orientations and phases

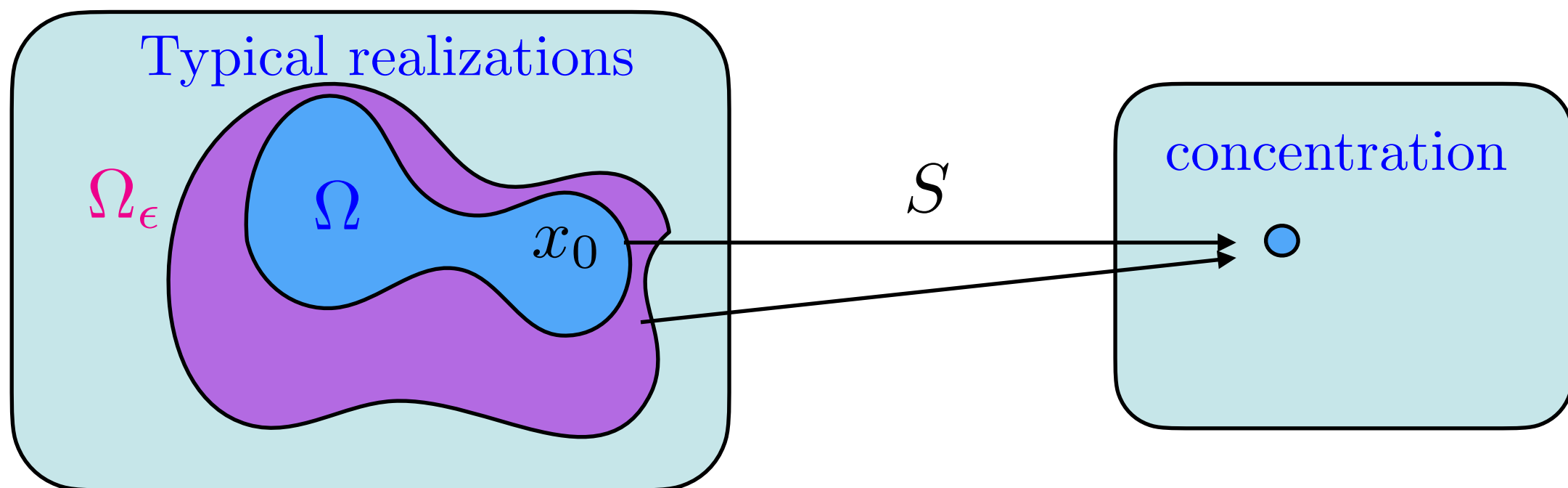**II- Image classification by deep separation and concentration:**

- Deep nets from priors without learning

- Learning tight frame contractions along channels only

Vector of statistics $S(x)$: observable

Concentration: $\text{Prob}_p \Big( \|S(x) - \mathbb{E}_p(S(x))\| > \epsilon \Big) \underset{d \to \infty}{\longrightarrow} 0$

$\Rightarrow$ a realisation $x_0$ satisfies $S(x_0) \approx \mathbb{E}_p(S(x))$ with high proba.



*Microcanonial ensemble:* $\Omega_\epsilon = \{x \ : \ \|S(x) - S(x_0)\| \leq \epsilon\}$

Maximum entropy model $\tilde{p}$ supported in $\Omega_\epsilon$ is uniform.

Generation by sampling $\tilde{p}$: SGD on $\|S(x) - S(x_0)\|$ from white noise
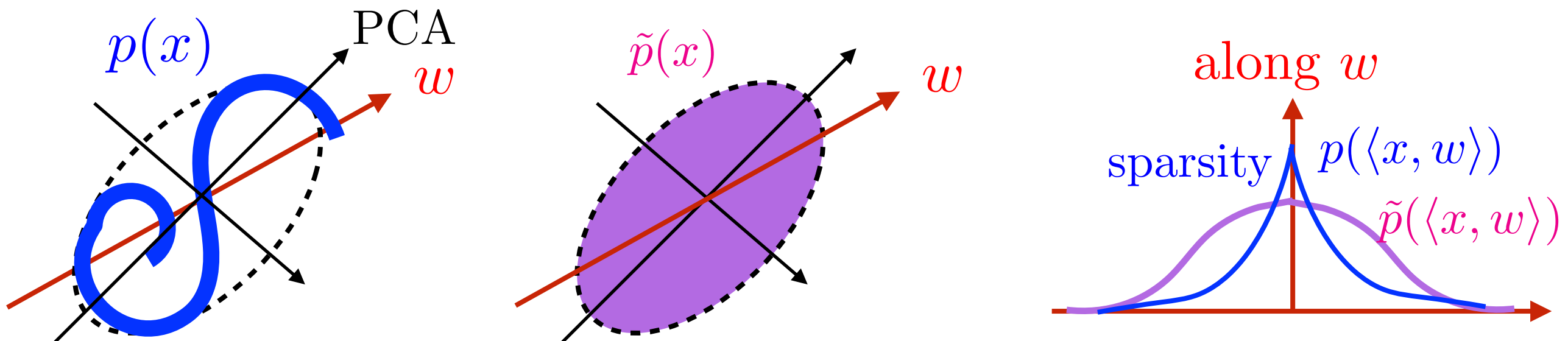not exactly maximum entropy *(J. Bruna)*

"Sufficient statistics" if $\Omega \approx \Omega_\epsilon$: how to define $S$ ?

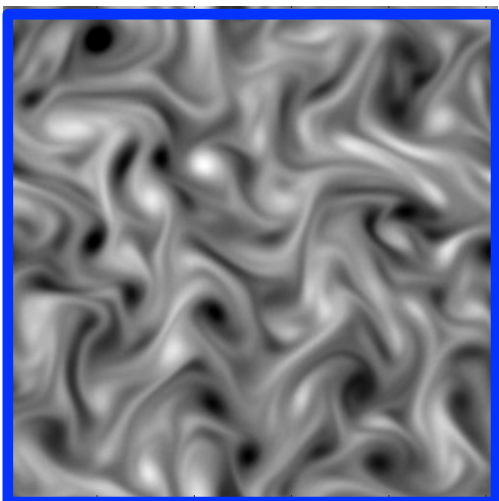Symmetry prior: $p(x)$ is translation invariant

$$S(x) = \left( d^{-1} \sum_u x(u)\, x(u - \tau) \right)_\tau$$

empirical covariance concentrates by spatial averaging

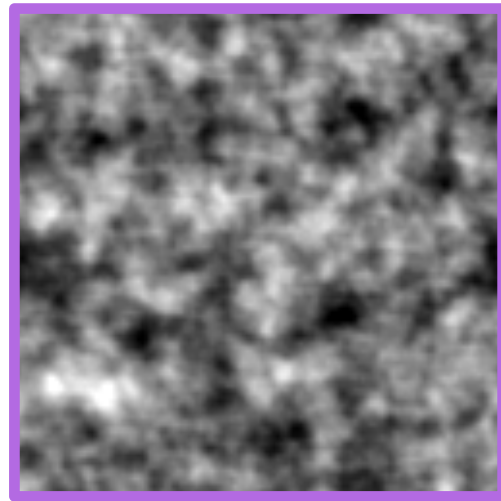Maximum entropy model $\tilde{p}$ asymptotically Gaussian: how good ?

PCA

$w$

along $w$
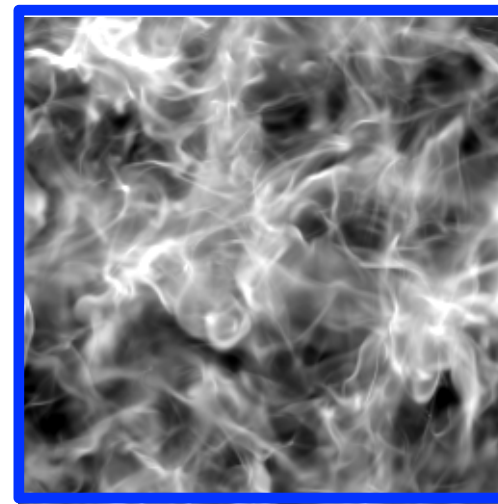
sparsity

$p(\langle x, w \rangle)$

$\tilde{p}(\langle x, w \rangle)$
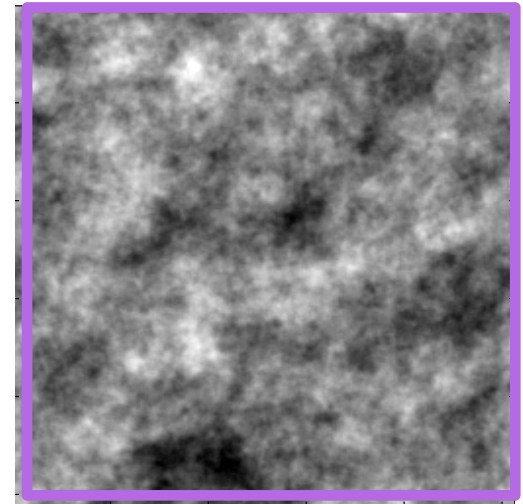
observation

model

Analysis

Fluide

Gaussien

Gaz

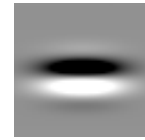Gaussien

# **Separation with Wavelets**

- Wavelet filter $\psi^\alpha(u)$:

  phase $\alpha$:    $0$        $\pi/2$

Scales $2^j$, angles $\theta$, phases $\alpha$: $\psi_\lambda(u) = 2^{-2j}\,\psi^\alpha(2^{-j} r_\theta u)$



Wavelets $\psi_\lambda$

- Wavelet tight frame separation: $Wx(u, \lambda) = x \star \psi_\lambda(u)$

- Not correlated across "channels" if $x$ is stationary:

$$\mathbb{E}\Big(Wx(u, \lambda)\,Wx(u, \lambda')\Big) \approx 0 \;\; \text{if} \;\; \lambda \neq \lambda'$$

Filters separation $F_w$

Tight frame
$F_w^T F = Id$

$2^0$

$F_w$

$F_w$

Sparse but not correlated
because of sign variations

$2^J$
Scale

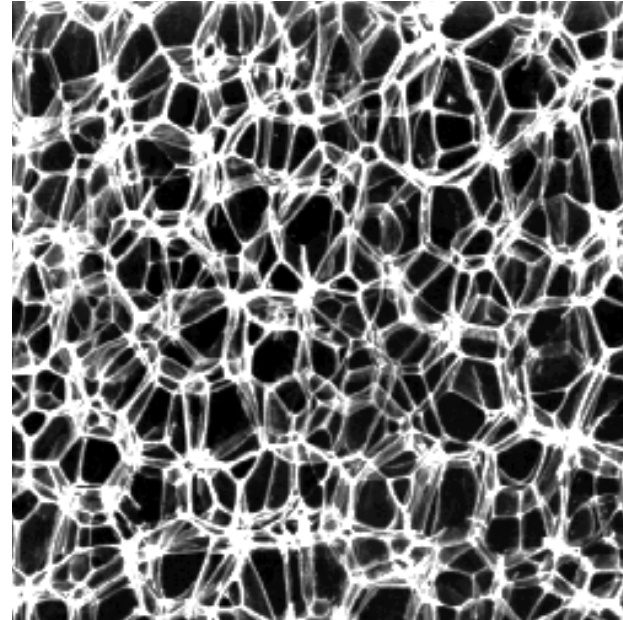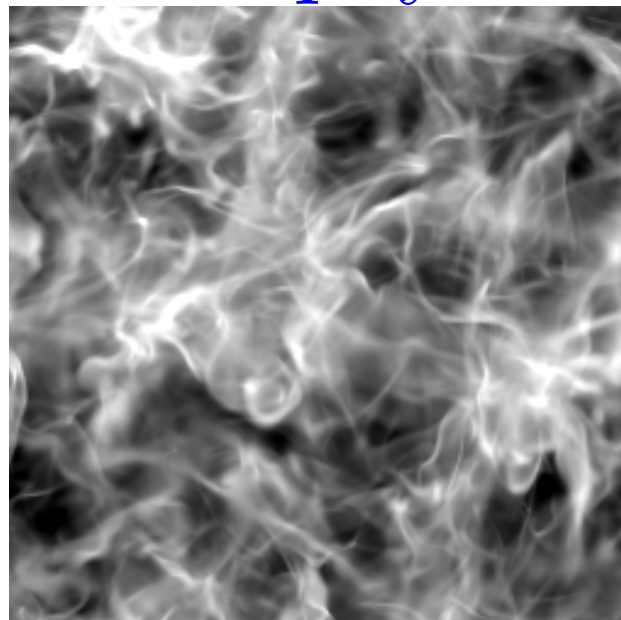How to capture dependance across scales, angles, phases channels ?

# Models of Stationary Processes

*Sixin Zhang*

$x_0$      Astrophysics      Ising-critical



Correlations across scales/orientations/phases $\lambda = (2^j, \theta, \alpha)$

are created by a ReLu $\rho(a) = \max(a, 0)$ (no bias) which shrinks sign:

$$S(x) = d^{-1} \sum_u \rho W x(u) \, \rho W x(u)^T \quad : \text{ empirical correlation}$$

Concentration by spatial averaging: dimension $O(\log^2 d)$

Maximum entropy models conditioned by $S(x_0)$

# Sampling from Max Entropy Model

*Sixin Zhang*
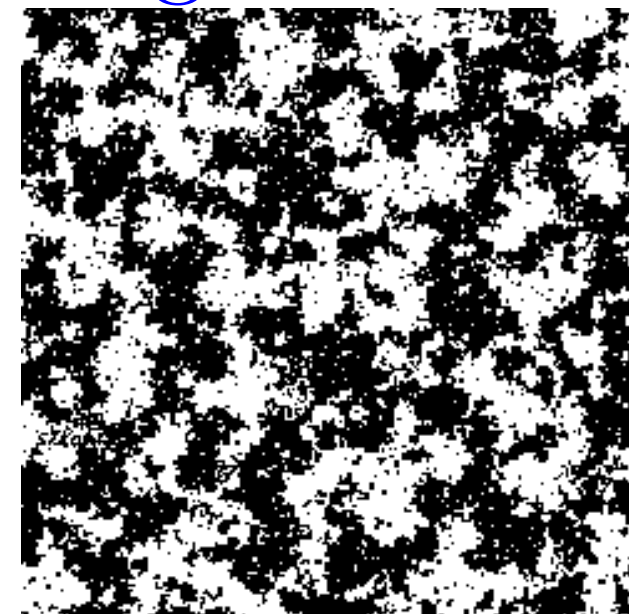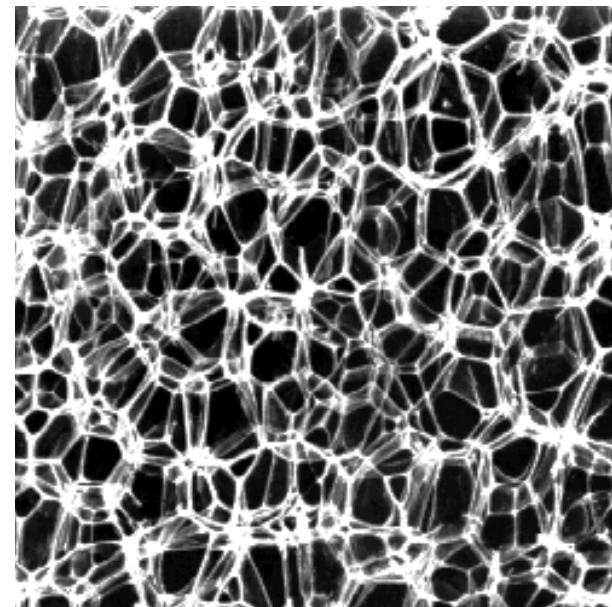
$d = 6\,10^4$

<span style="color:blue">Astrophysics</span>    <span style="color:blue">Ising-critical</span>

$x_0$

$x$



<span style="color:red">$S(x_0)$ has $2\,10^3$ empirical covariances</span>

Sampled from $S(x_0)$ with $SGD$ algorithm

*E. Allys, T. Marchand, J.F. Cardoso, F. Villaescusa, S. Ho, S. Mallat*

Generation of matter density fields from rectified wavelet covariances:



Original $x_0$ — Max-entropy generation

- Reproduces high order moments

- Accurate regression of 6 cosmological parameters from $S(x_0)$

- A deep network progressively separates and concentrates

  - Can we do it from prior without learning ?

  - If not, what needs to be learned ?

Wavelet separation:

$$\rho\, W_1 x \equiv \begin{pmatrix} x \star \phi_{2^J} \\ \rho(x \star \psi_{\lambda_1}) \end{pmatrix}_{\lambda_1}$$

Concentration by averaging

$x(t)$

$x \star \phi_{2^J}(t)$

$2^J$

Lost high frequencies: $x \star \psi_{\lambda_1}(t)$ but $(x \star \psi_{\lambda_1}) \star \phi_J = 0$

Relu non-linearity: $\rho(x \star \psi_{\lambda_1}(t))$ to remove sign

Concentration: $\rho(x \star \psi_{\lambda_1}) \star \phi_{2^J}(t)$

To preserve separation

Need to recover lost high frequencies: $\rho(x \star \psi_{\lambda_1}) \star \psi_{\lambda_2}(t)$

Concentration with Relu and averaging:

$$\begin{pmatrix} \rho(x \star \psi_{\lambda_1}) \star \phi_{2^J}(t) \\ \rho(\rho(x \star \psi_{\lambda_1}) \star \psi_{\lambda_2}) \star \phi_{2^J}(t) \end{pmatrix}_{\lambda_2}$$

# Wavelet Scattering Network



Frame contraction: $\rho\, F_w$

$F_w^T F = Id$

$\rho\, F_w$

$\rho(x \star \psi_{\lambda_1})$

$\rho\, F_w$

$\rho(\rho(x \star \psi_{\lambda_1}) \star \psi_{\lambda_2})$

$\rho\, F_w$

$\rho\, F_w$

Depth $= J$

channels

$\Phi = (\rho\, F_w)^J$ : iterated frame contractions

Scatters along progressively more channels

A convolution tree: no channel interactions, no learning.

# Image Classification

Scat-Net$_J$ : $J$ layers

$$x \rightarrow \boxed{\rho(\rho(x \star \psi_\lambda) \star \psi_{\lambda'}) \star \phi_J} \rightarrow \Phi(x) \rightarrow \boxed{C} \rightarrow \tilde{y}$$

no learning

supervised learning

channels     $2^J$

Errors:

| | Scattering | Deep Nets. |
|---|---|---|
| MNIST: $28^2$<br>10 classes | $J = 3$   0.5 % | 0.5 % |
| CIFAR: $32^2$<br>10 classes | $J = 4$   23% | ResNet-18: 8%<br>ResNet-50: 7.6% |
| ImageNet: $228^2$<br>$10^3$ classes<br>1 million training | 52 %<br>$J = 6$ | AlexNet-7: 20%<br>ResNet-18: 11%<br>Res-Net 50: 7% |

mite     container ship     motor scooter     leopard

grille     mushroom     cherry     Madagascar cat

What is learned ?

*John Zarka, Florentin Guth*

Frame soft-thresholding along scattering channels:

$$C_1 C_1^T = Id \qquad F_1^T F_1 = Id$$

orthog. 1x1 conv.     tight frame

spatial wavelet scattering     concentration     soft-thresholding

$$x \rightarrow \boxed{\text{Scat-Net}_J} \xrightarrow{\Phi(x)} \boxed{C_1} \rightarrow \boxed{F_1^T \rho F_1} \rightarrow \boxed{C} \rightarrow \tilde{y}$$

1200     256     256

- SGD optimisation

| $\Phi(x)$ | | Scat. | 1CoScat | ResNet-18 |
|---|---|---|---|---|
| CIFAR | Error | 27% | 18% | 8% |
| | Fisher | 22 | 30 | |
| ImageNet Top 5 | Error | 60% | 30% | 11% |
| | Fisher | 2.9 | 3.4 | |

# Concentrated Scattering

*J. Zarka, F. Guth*

$x \in \mathbb{R}^d$ $\quad$ $F_1^T \rho F_1 C_1$ $\quad$ $F_2^T \rho F_2 C_2$ $\quad$ $F_J^T \rho F_J C_2$

$\rho F_w$ $\quad$ $\rho F_w$ $\quad$ $\rho F_w$ .... $\rho F_w$ $\quad$ $C$ $\quad$ $max$ $\tilde{y}$

3 $\qquad$ 64 channels $\qquad$ 128 channels $\qquad$ ... 512 ...

- Network without learning bias and semi-orthogonal operators

- Learning 1x1 convolutions across scattering channels

- SGD optimisation

| $\Phi(x)$ | | 1CoScat | CoScat | ResNet-18 |
|---|---|---|---|---|
| CIFAR | Error | 18% | 7.8% | 8% |
| | Fisher | 30 | 70 | |
| | Depth | 5 | 8 | 18 |
| ImageNet Top 5 | Error | 30% | 13% | 11% |
| | Fisher | 3.4 | 7.2 | |
| | Depth | 7 | 12 | 18 |

Mathematical control of Fisher ratios ?

# Conclusion

- **Deep network separate and concentrate: what mechanism ?**

- Variance can be reduced with tight frame contractions

- Spatial filtering can be handled with wavelet frame which separate scale, angle and phase channels.

- Learning contractions along channels can reach ResNet accuracy

- Control of *Fisher ratios* is an open math. problem.

New Interpretable Statistics for Large Scale Structure Analysis and Generation

*Allys, Marchand, Cardoso, Villaescusa, Ho, Mallat,* arXiv:2006.06298, Phys. Rev.

Tight Frame Contractions in Deep Networks

*J. Zarka, F. Guth,, S. Mallat,* OpenReview, ICLR 2021