# A Mathematical Quest

# for Information Processing

*Stéphane Mallat*

**Collège de France**

COLLÈGE DE FRANCE —1530—

# A Century of Information

What is information in mathematics ?

1921
Fisher Information
Math. Statistics

1948
Shannon Information Theory
Proba. Concentration

How to represent data and analyse information ?

1940's
Wiener Gaussian Models
Fourier Analysis

1990's
Wavelet Sparse Models
Functional Analysis

2021
Neural networks
High-dim. Geometry
Maths not understood

1921 report of Ronald Fisher to the Royal Academy

"The object of statistical methods is the reduction of data"

"Represent the whole of the relevant data information"

- Model data $\{x_t\}_{t \le n}$ as independant samples
  of a distribution $p_\theta(x_t)$ parameterised by $\theta$.
  Gaussian example: $p_\theta(x) \sim e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ with $\theta = (\mu, \sigma)$.

- Consistent estimator $\hat{\theta} \to \theta$ as $n$ tends to $\infty$

- Maximum likelihood estimator $\hat{\theta}$ of $\theta$ given $\{x_t\}_{t \le n}$

  $$\hat{\theta} \text{ maximises } p_\theta(x_1, ..., x_n) = \prod_t p_\theta(x_t)$$

- Amount of information carried by data with probability $p_\theta$ on the unknown parameter $\theta$: curvature of $\log p_\theta$

$$I(\theta) = \mathbb{E}\left[\left(\frac{\partial \log p_\theta(x)}{\partial \theta}\right)^2\right]$$

- **Cramer-Rao Bound** on parameter estimation (1940's):

**Theorem** If $\mathbb{E}(\hat{\theta}) = \theta$ then
$$\mathbb{E}(\hat{\theta} - \theta)^2 \geq \frac{1}{I(\theta)}$$

The Fisher information controls the uncertainty to estimate $\theta$

What family of parametrised probabilities $\{p_\theta\}_\theta$ ?

*Concentration in high dimension*

$n$ independent random variables $X = (X_1, ..., X_n)$

with same probability distribution $p(X) = \prod_t p(X_t)$

Law of large numbers:        *Entropy*

$$-\frac{1}{n} \log p(X) = -\frac{1}{n} \sum_{t=1}^{n} \log p(X_t) \xrightarrow[n \to \infty]{} H = \mathbb{E}\Big( -\log p(X_t) \Big)$$

## Asymptotic Equipartion Theorem

For an ergodic stationary process $\{X_t\}_t$

$$-\frac{1}{n} \log p(X_1, ..., X_n) \xrightarrow[n \to \infty]{} H \quad \text{with probability } 1$$

Typical set: $T_\epsilon = \left\{ x \in \mathbb{R}^n \; : \; |n^{-1} \log p(x) - H| \leq \epsilon \right\}$

*Concentration:* $\mathrm{Prob}(X \in T_\epsilon) \xrightarrow[n \to \infty]{} 1$

If $x \in T_\epsilon$ then $p(x) \sim 2^{-nH}$
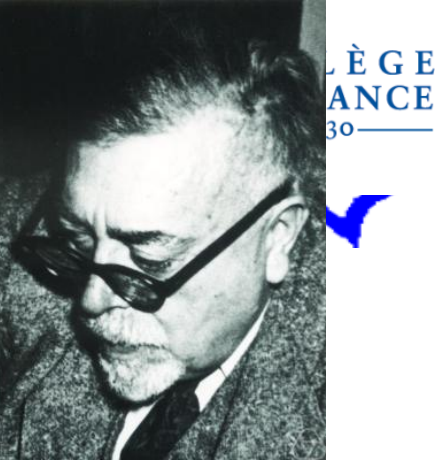
If $X$ is quantised then $|T_\epsilon| \sim 2^{nH}$

$\mathbb{R}^n$

$T_\epsilon$

$nH$ is the minimum average number of bits to code $X$

## Considerable impact:

- Coding: telecommunication and data storage

- Statistical physics (thermodynamic entropy)

- Large Deviation Theory (*Donsker-Varadhan 1960's*)

How to specify the geometry of Typical sets ?

$$p_\Theta(x) = Z^{-1} \exp\left(-\tfrac{1}{2}\langle x, \Theta x\rangle\right)$$

*Wiener*

where $\Theta$ is a positive matrix of parameters

Typical sets $T_\epsilon$ are ellipsoids whose principal axes are vectors vectors of an orthonormal basis $\mathcal{B}$ which diagonalises $\Theta$.



- If $X_1\dots,X_t,\dots$ is stationary, i.e. $p(x)$ is invariant to time-shift then $\mathcal{B}$ is a Fourier basis: $X_t = \sum_\omega \tilde{X}_\omega\, e^{it\omega}$

Limit of continuous time: spectral representation
Typical sets: balls of weighted Fourier spaces (Sobolev).

COLLÈGE DE FRANCE 1530

- Non Gau

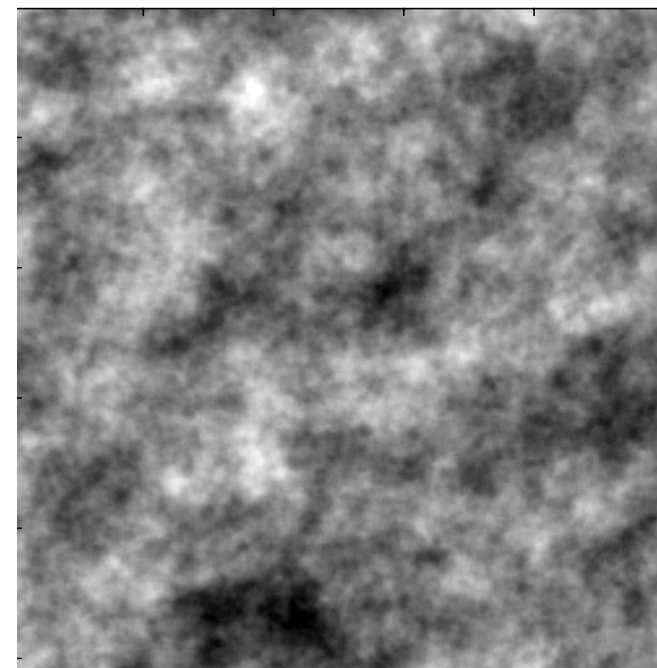Original          Gaussian model (Kolmogorov)
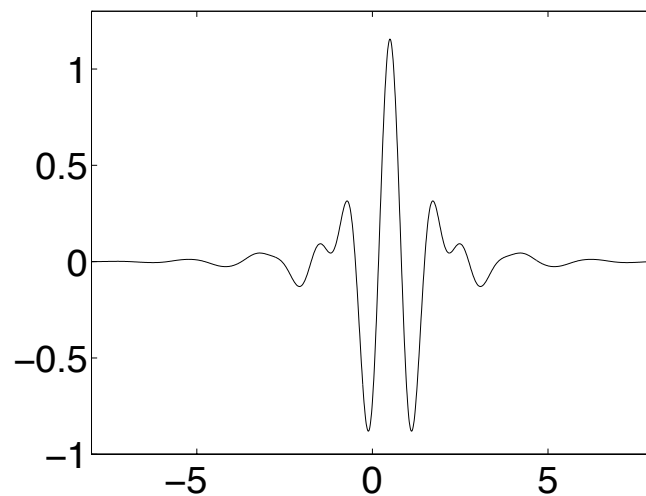
(1941)

Fluid
Turbulence

Same power
spectrum

Cosmologic
Turbulence

How to represent X and identify typical sets ?

- Represent transient phenomena with localised wavelets.

- Sparse representations in wavelet bases (1980-90's):

Meyer wavelet $\psi(t)$

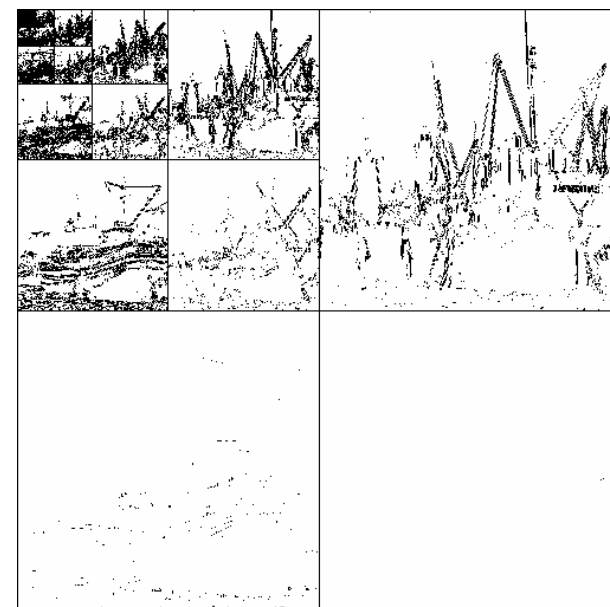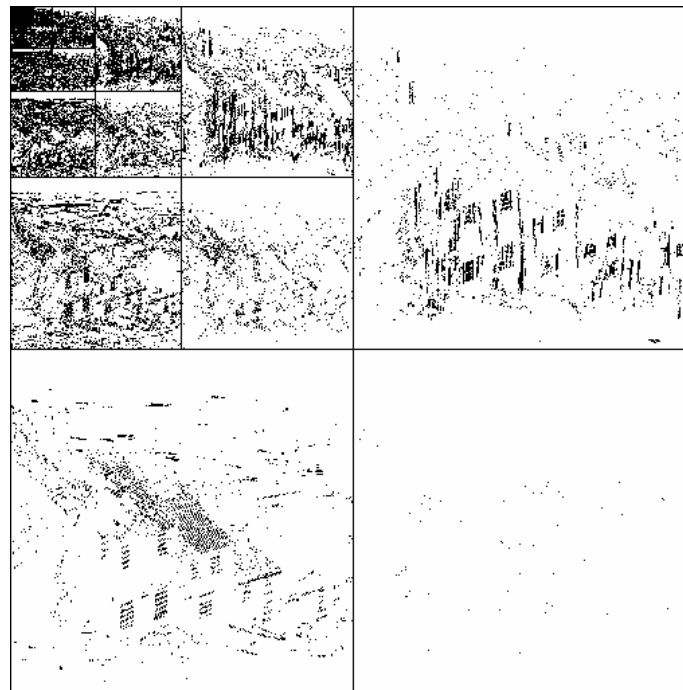Orthonormal basis of $\mathbf{L^2}(\mathbb{R})$
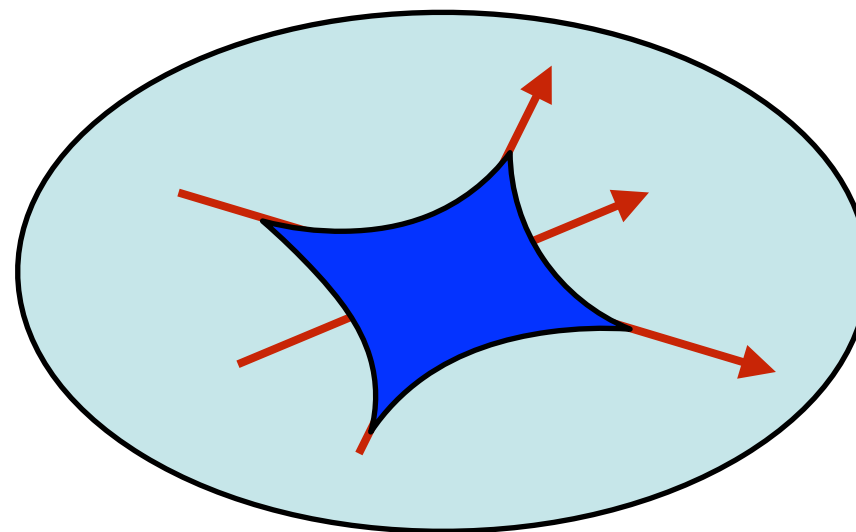


$X$

Sparse wavelet coefficients

Original

Sparse
Wavelet coefficients

JPEG-2000
Compressed by 40

Typical sets $T_\epsilon$ : bounded weighted $\ell^p$ norms of wavelet coef.
balls of Besov spaces

Still too crude to model geometric image structures: what else ?
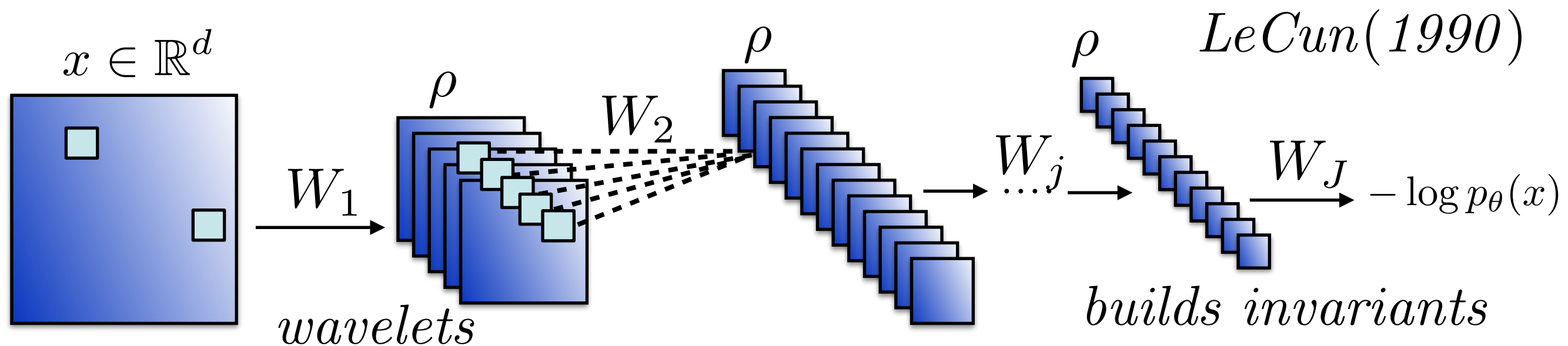
*McCulloch and Pitts (1943)*

- Alternate linear operators and a pointwise non-linearity:

$$-\log p_\theta(x) = W_J \, \rho \, W_{J-1} \, ... \, \rho W_2 \, \rho \, W_1 x$$
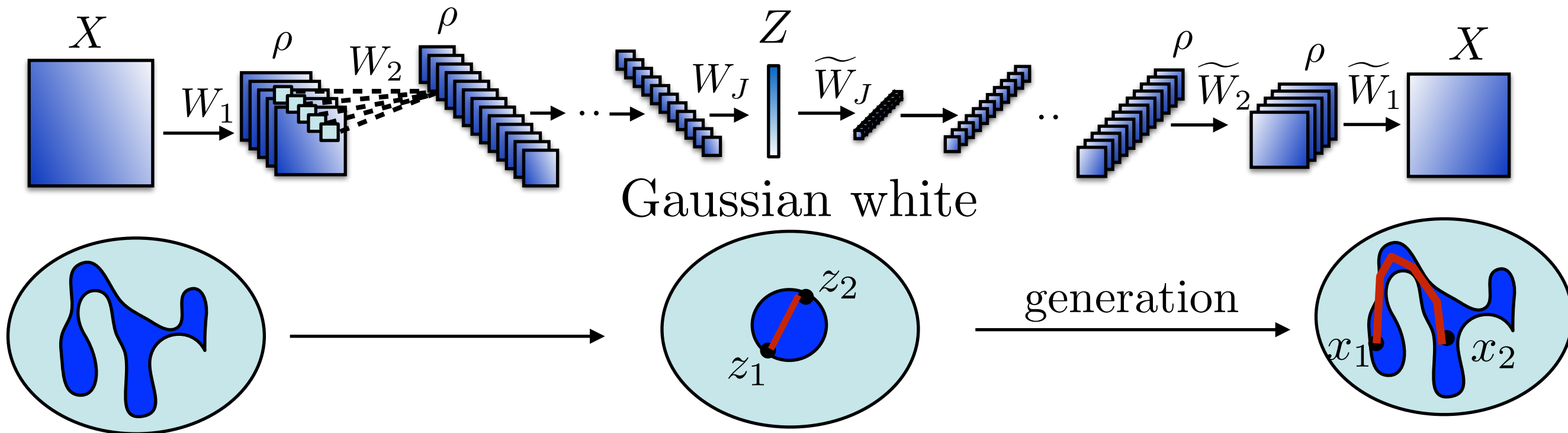
with a rectifier $\rho(\alpha) = \max(\alpha, 0)$

and $\theta = (W_j)_{1 \leq j \leq J}$ are matrices optimised by

maximising the data *likelihood* with a gradient descent.

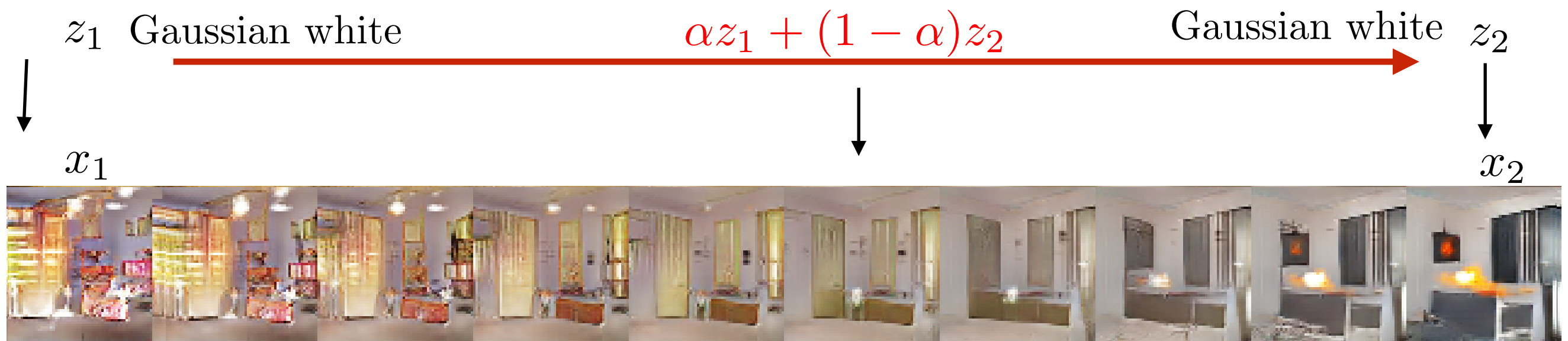- Convolutional architectures: shift-invariant operators $W_j$



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

$\rho$    $W_2$    $\rho$    $\rho$    *LeCun(1990)*

$W_1$    $\underset{....}{W_j}$    $\underset{}{W_J}$    $-\log p_\theta(x)$

*wavelets*          *builds invariants*

mapping into a Gaussian and inversion



Gaussian white



generation

Can generate complex ergodic processes including turbulences
Beyond ergodicity: generation from images of bedrooms

$z_1$ Gaussian white          $\alpha z_1 + (1 - \alpha) z_2$          Gaussian white $z_2$

$x_1$                                                                          $x_2$
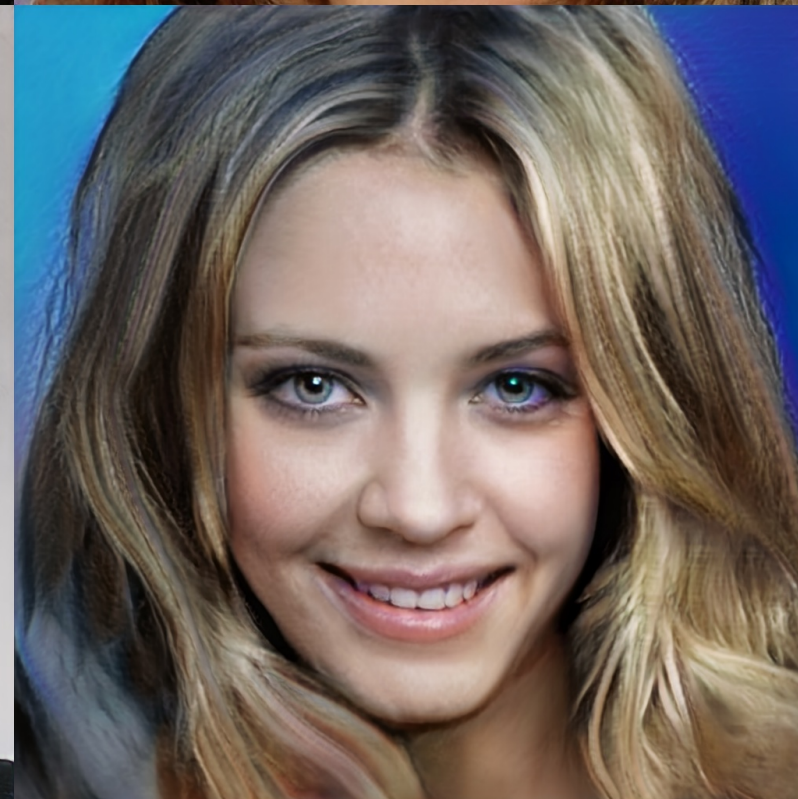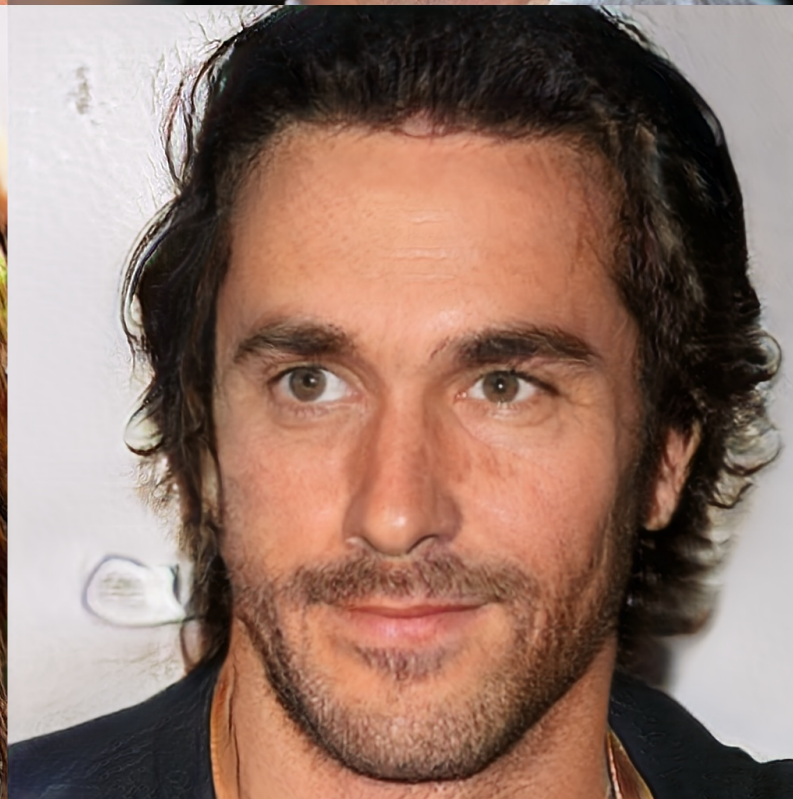


What class of random processes / transport / function spaces ?
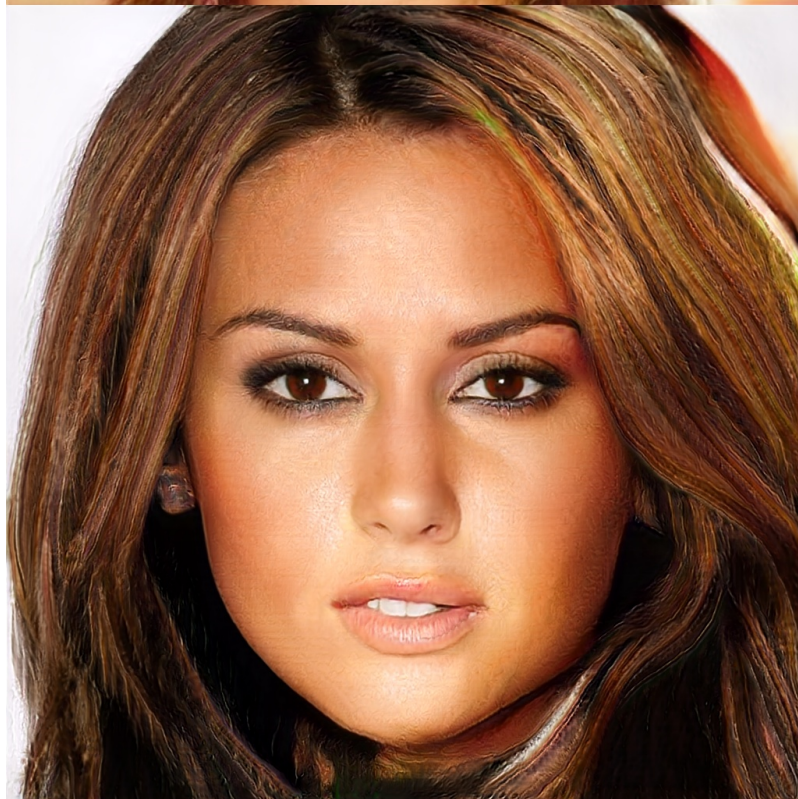
# High Resolution Generation

*T. Karras, T. Aila, S. Laine, J. Lehtinen*

Generated from Hollywood celebrities data basis

# Outstanding Mathematical Questions

- Information processing is about high-dimensional geometry.

- Neural networks have spectacular ability to process information, but mathematically not understood.

- Major societal issue because of critical AI applications: *medical, transport, decision making...*

- Outstanding questions, from *statistics* to:
  - Probability and concentration
  - Functional and harmonic analysis
  - Geometry and group theory
  - Optimisation and high-dimensional transport