AIMS CDT - Signal Processing Michaelmas Term 2025

Xiaowen Dong

Department of Engineering Science



Logistics

- Lectures
 - Monday-Thursday 10.00-12.30 @IEB
 - slides: http://www.robots.ox.ac.uk/~xdong/teaching.html
- Lab sessions
 - Monday-Thursday 14.00-17.00 @EH
 - notes: http://www.robots.ox.ac.uk/~xdong/teaching.html
 - demonstrators
 - Yin-Cong Zhi (yin-cong.zhi@ndph.ox.ac.uk)
 - Scott le Reux (scott.leroux@wolfson.ox.ac.uk)
 - Ning Zhang (ning.zhang@some.ox.ac.uk)
 - Jacob Bamberger (jacob.bamberger@some.ox.ac.uk)
- Questions & Comments: xdong@robots.ox.ac.uk

Overview

- Lecture 1: Introduction to signal processing
 - time-frequency analysis, filtering, Fourier & wavelet transforms, dictionary learning
 - Lab 1: Signal and image processing
- Lecture 2: Introduction to graph signal processing
 - graph Fourier transform, filtering & convolution, representation of graph signals
 - Lab 2: Graph signal processing
- Lecture 3: Deep learning on graphs
 - convolutional neural networks on graphs, message passing neural networks
 - Lab 3: Graph neural networks
- Lecture 4: Bayesian modelling of graph-structured data
 - Gaussian processes on graphs, Bayesian optimisation of graph-based functions
 - Lab 4: Gaussian processes on graphs

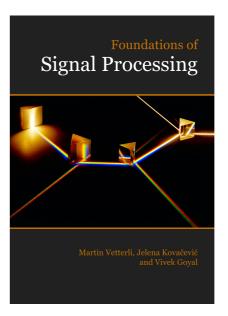
Resources

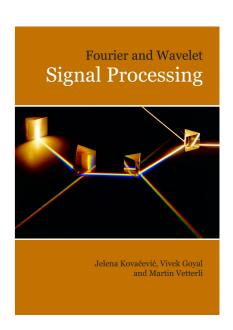
Textbooks

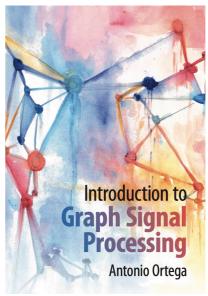
- Vetterli et al. Foundations of signal processing.
 Cambridge University Press, 2014. Available at http://www.fourierandwavelets.org
- Kovačević et al. Fourier and wavelet signal processing. Available at http://www.fourierandwavelets.org
- Ortega. Introduction to graph signal processing.
 Cambridge University Press, 2022.
- Hamilton. Graph representation learning. Morgan & Claypool Publishers, 2020. Available at https://www.cs.mcgill.ca/~wlh/grl_book/

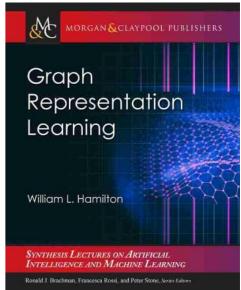
Resources

- https://web.media.mit.edu/~xdong/resource.html
- https://github.com/naganandy/graph-based-deep-learning-literature









Introduction to Signal Processing

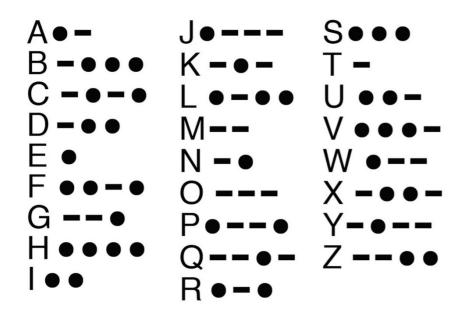
Lecture 1

- Introduction & Basic concepts and tools
- A historical overview of signal representation techniques
- Applications & Discussion

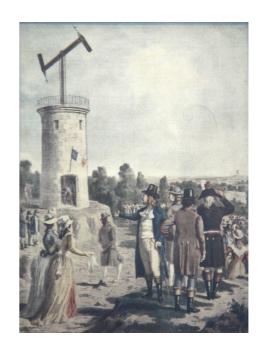
Historical notes



smoke signal tower (1570)



Morse code (1830s)

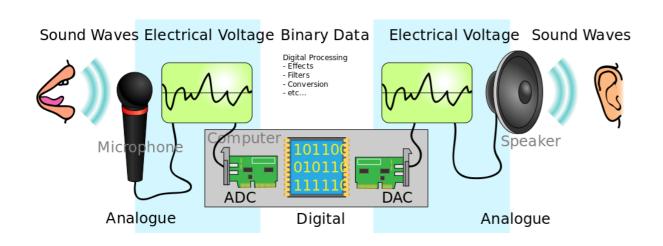


semaphore telegraph (1792)

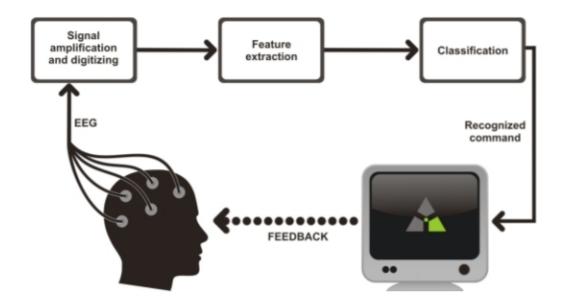


electronic communication (today)

Modern signal processing applications



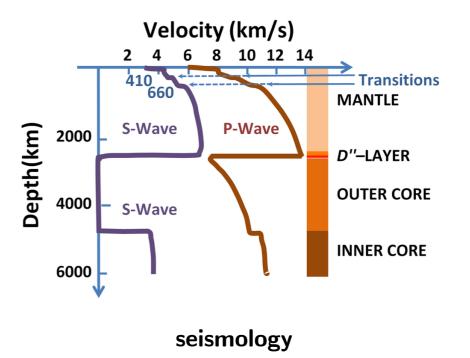
speech processing



EEG signal classification



image denoising



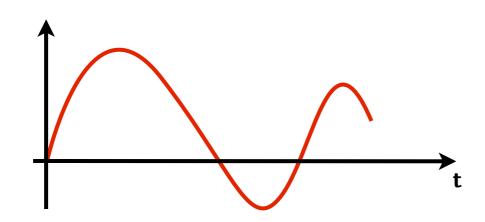
discrete

continuous

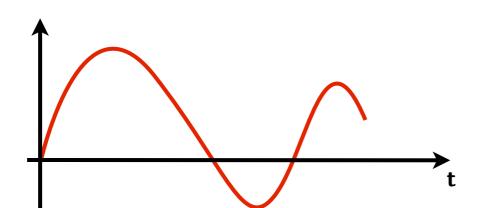
time

discrete

continuous

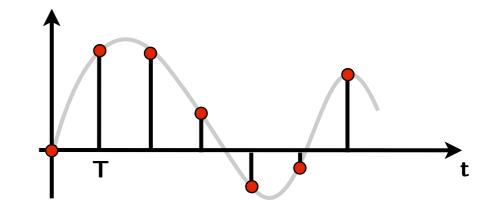


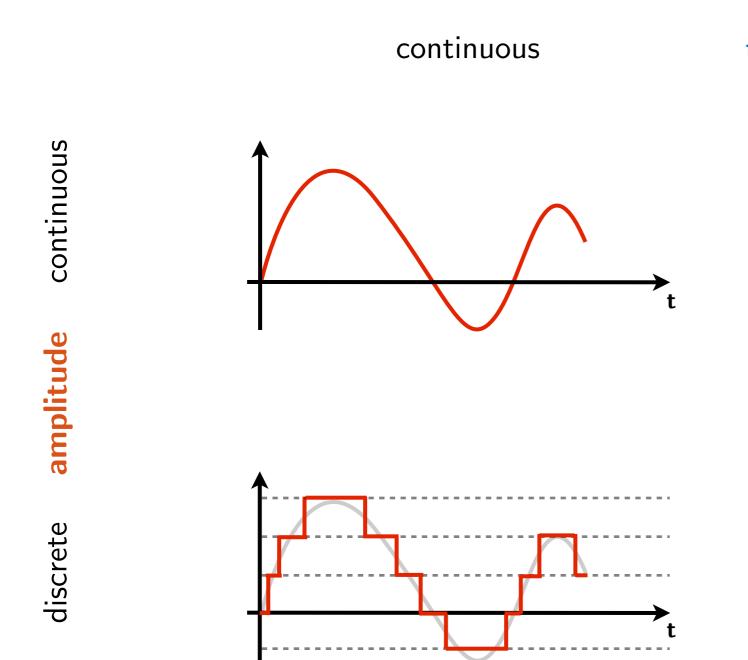
continuous

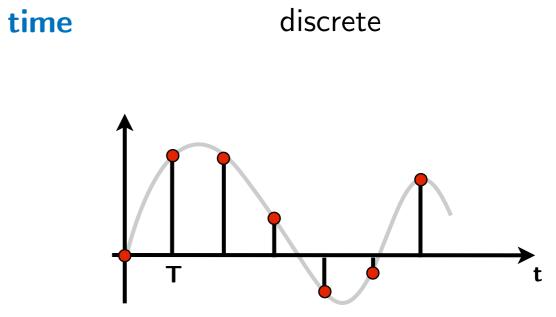


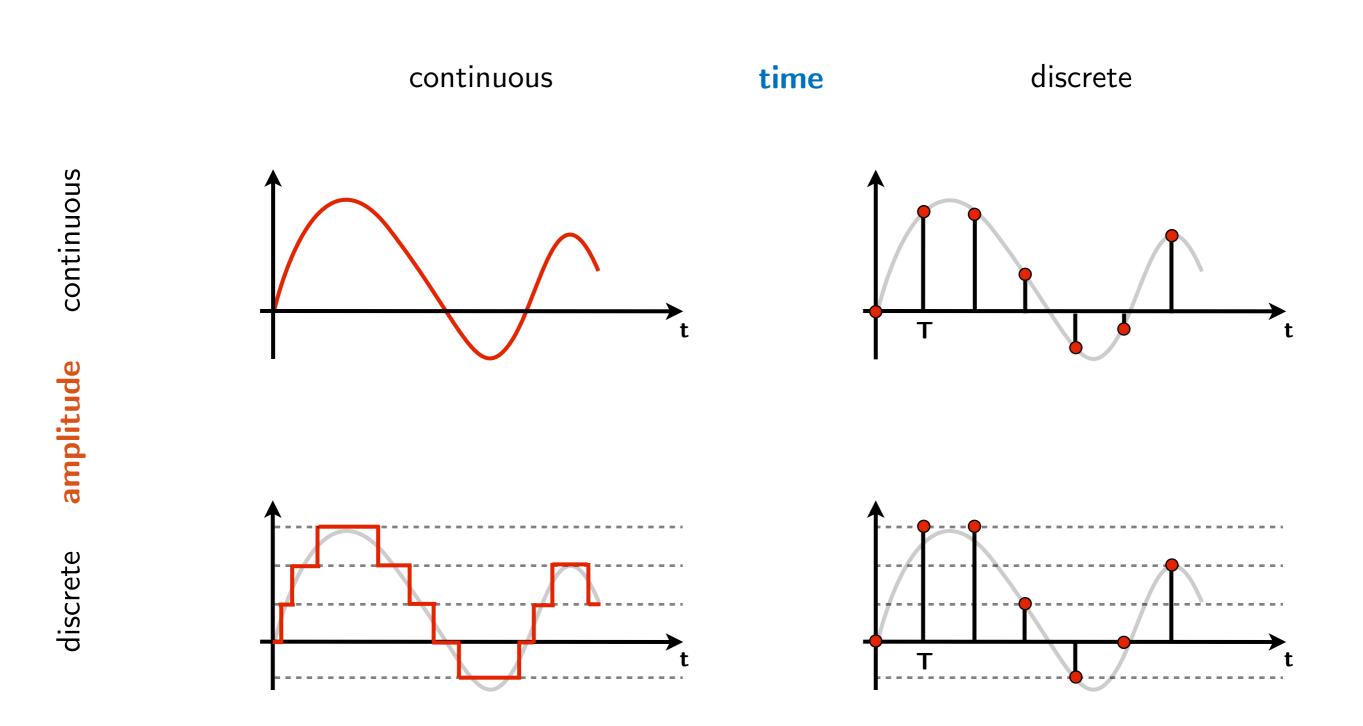
time

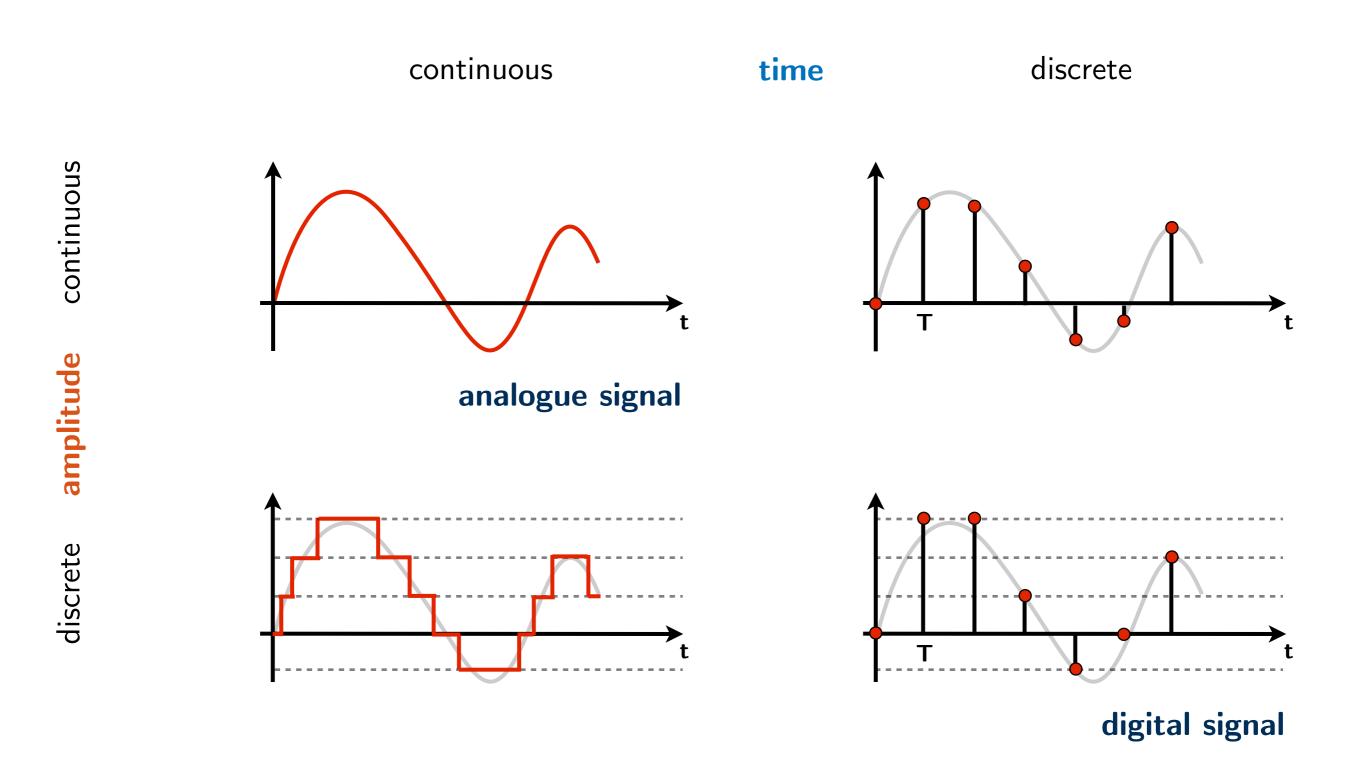
discrete

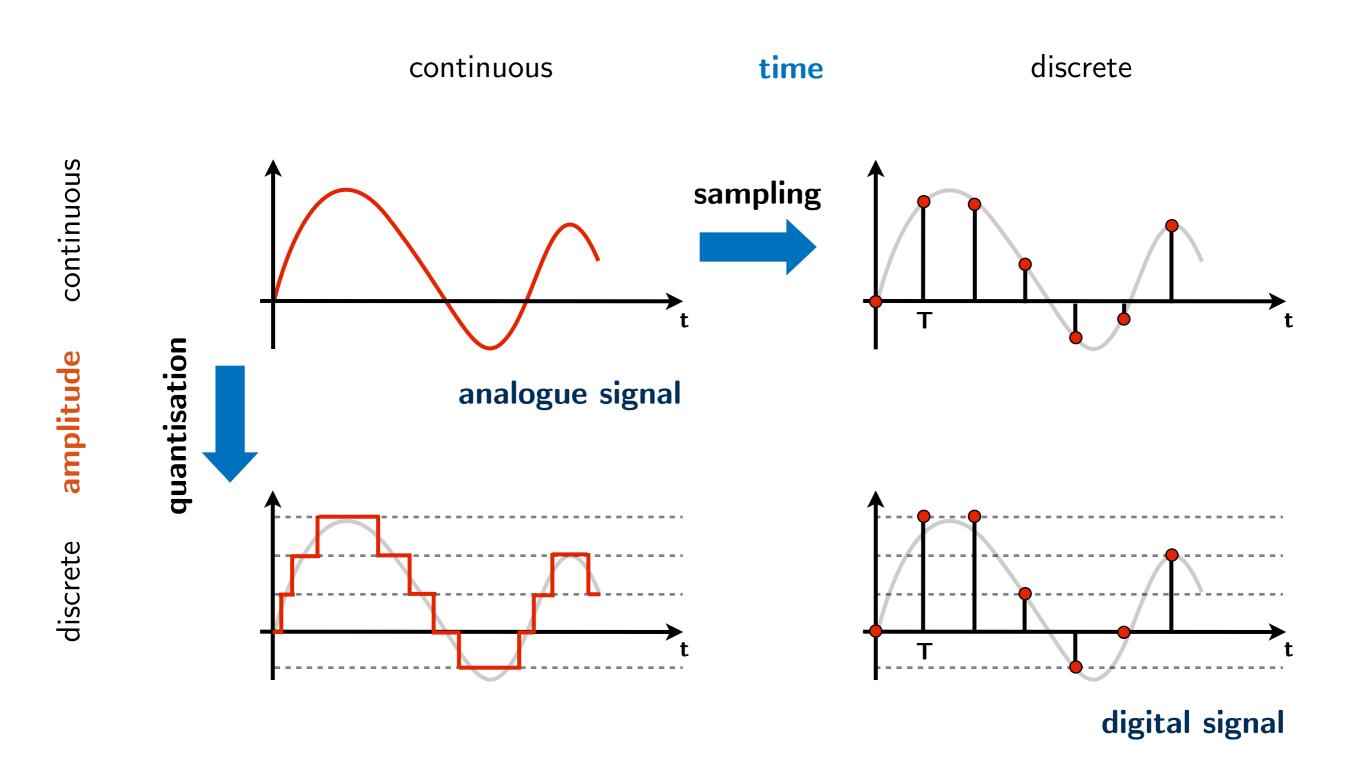






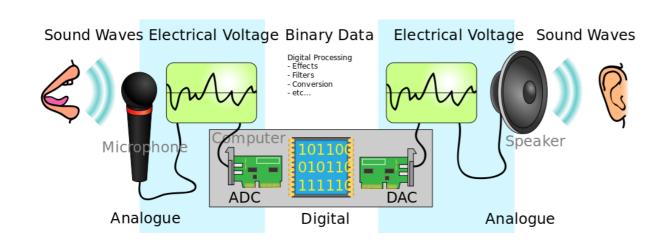


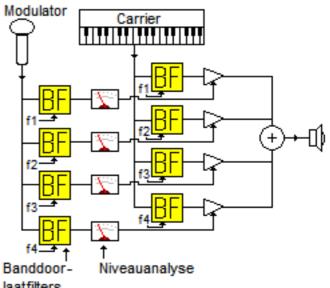




Analogue vs Digital signal processing

- Many signals of practical interest are analogue: e.g., speech, seismic, radar, and sonar signals
- Analogue signal processing systems are based on analogue equipment:
 e.g., channel vocoder
- Dramatic advance of digital computing moves the trend towards digital systems





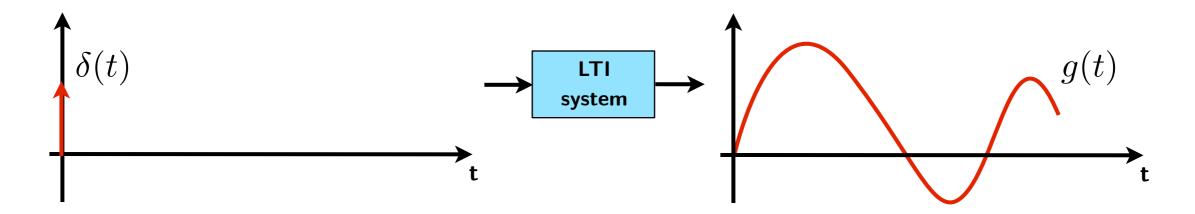
Signal processing as linear processes



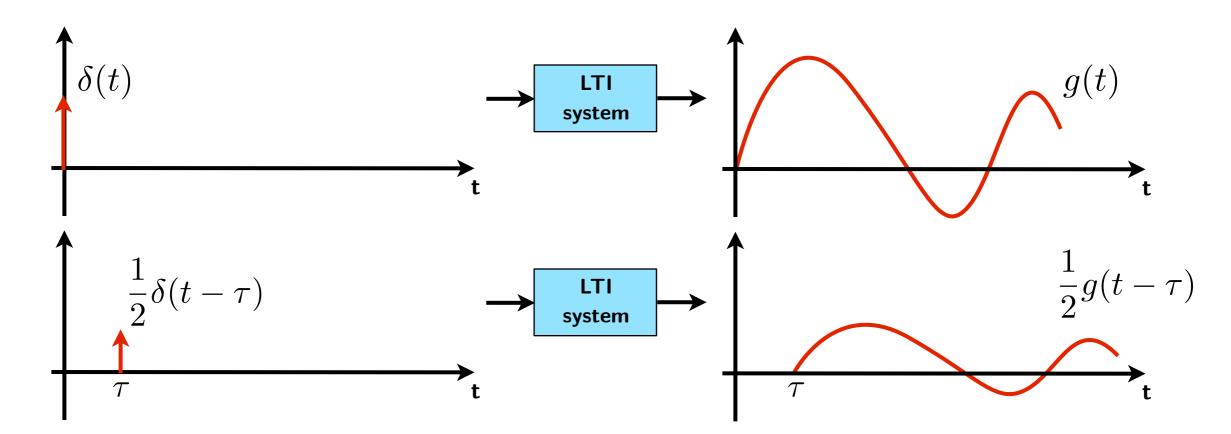
amplification/attenuation, filtering, (un-)mixing, etc.

- Linear time-invariant (LTI) system whose input-output characteristics can be defined by
 - impulse response in time domain
 - transfer function in frequency domain
- There is an invertible mapping between time- and frequency-domain representations

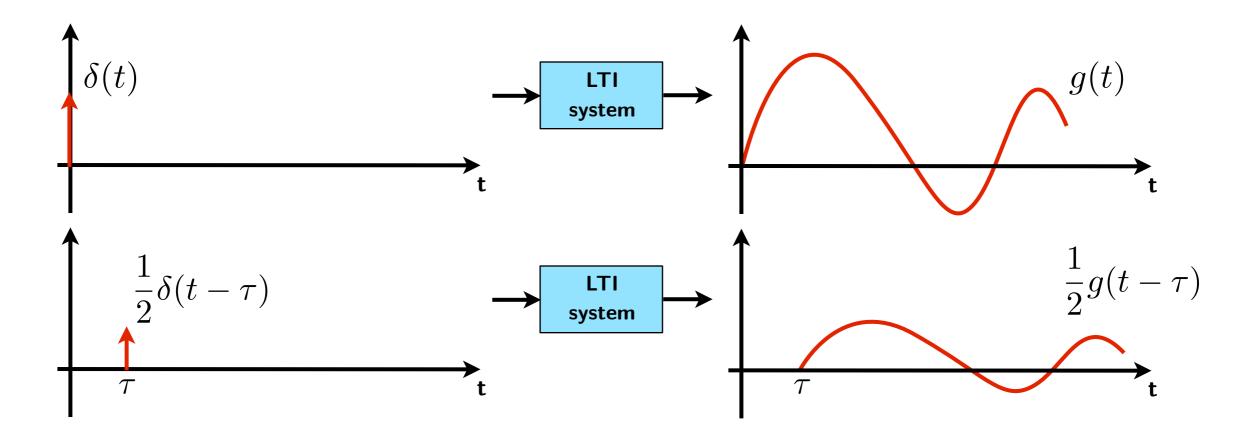
 Convolution allows the evaluation of the output signal from an linear time-invariant (LTI) system, given its impulse response and input signal



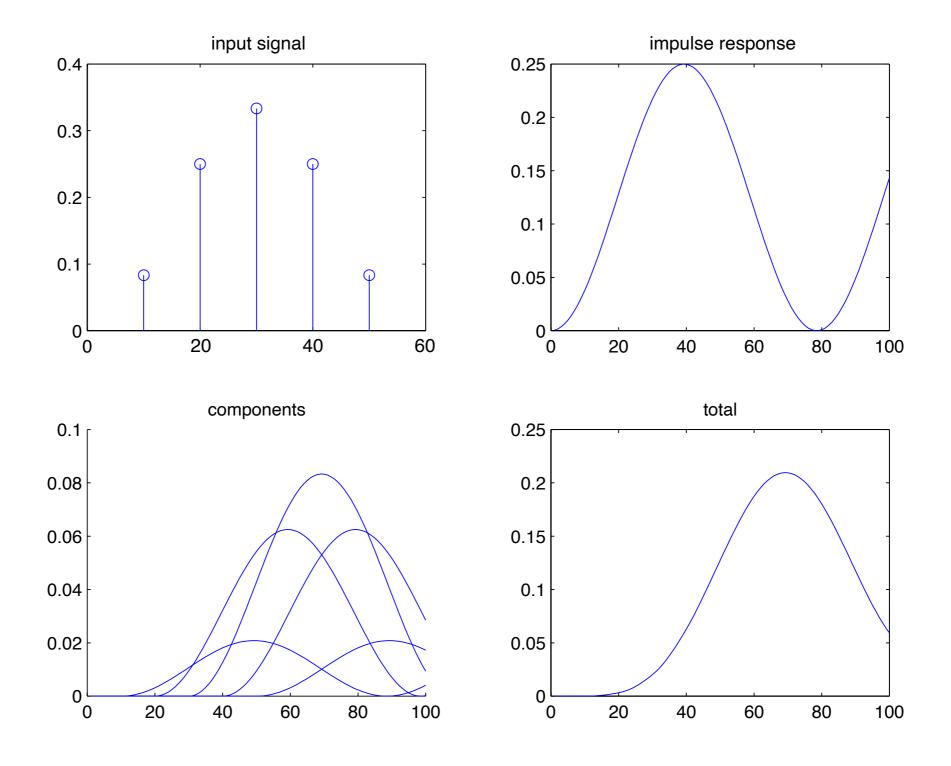
 Convolution allows the evaluation of the output signal from an linear time-invariant (LTI) system, given its impulse response and input signal

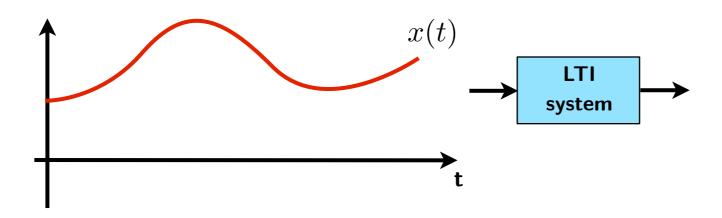


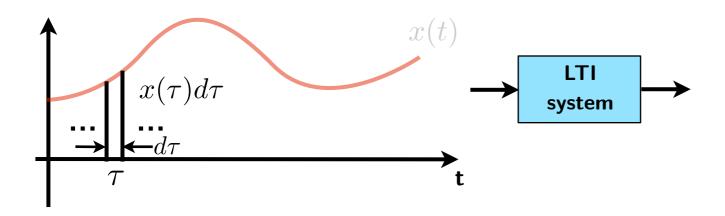
 Convolution allows the evaluation of the output signal from an linear time-invariant (LTI) system, given its impulse response and input signal

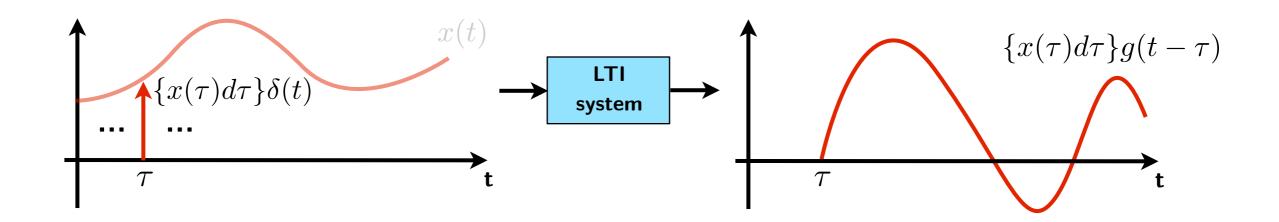


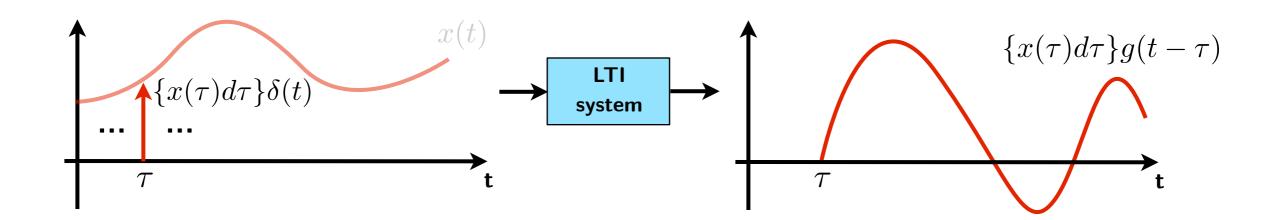
- Evaluate system output for
 - input: succession of impulse functions (which generate weighted impulse responses)
 - output: sum of the effect of each impulse function







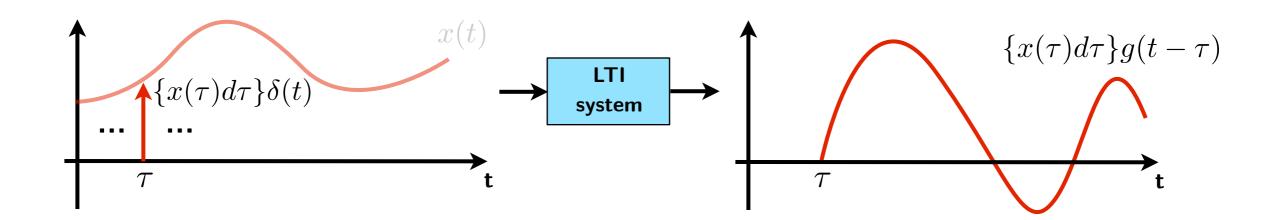




- this gives the convolution integral

$$y(t) = \sum_{\tau} \{x(\tau)d\tau\}g(t-\tau) \xrightarrow{d\tau \to 0} \int_{0}^{\infty} x(\tau)g(t-\tau)d\tau$$

- system output is convolution of input and impulse response



- this gives the convolution integral

$$y(t) = \sum_{\tau} \{x(\tau)d\tau\}g(t-\tau) \xrightarrow{d\tau \to 0} \int_{0}^{\infty} x(\tau)g(t-\tau)d\tau$$

- system output is **convolution** of input and impulse response
- impulse response characterises the system in time domain

Frequency-domain analysis

Consider the following LTI system

$$x(t) = e^{j\omega t} \longrightarrow g(t)$$

$$y(t) = \int_{-\infty}^{\infty} e^{j\omega(t-\tau)} g(\tau) d\tau = e^{j\omega t} G(j\omega)$$

- $e^{j\omega t}$ is an eigenfunction of an LTI system with eigenvalue $G(j\omega)$, which is the Fourier transform of the impulse response g(t)

Frequency-domain analysis

Consider the following LTI system

$$x(t) = e^{j\omega t} \longrightarrow g(t)$$

$$y(t) = \int_{-\infty}^{\infty} e^{j\omega(t-\tau)} g(\tau) d\tau = e^{j\omega t} G(j\omega)$$

- $e^{j\omega t}$ is an eigenfunction of an LTI system with eigenvalue $G(j\omega)$, which is the Fourier transform of the impulse response g(t)
- $G(j\omega)$, the frequency response, characterises the system in frequency domain

Fourier transform

• Fourier transform (FT) and inverse FT

$$X(j\omega) = \int_{-\infty}^{\infty} x(t)e^{-j\omega t}dt \qquad x(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} X(j\omega)e^{j\omega t}d\omega$$

Fourier transform

• Fourier transform (FT) and inverse FT

$$X(j\omega) = \int_{-\infty}^{\infty} x(t)e^{-j\omega t}dt \qquad x(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} X(j\omega)e^{j\omega t}d\omega$$

• FT of output is multiplication of FT of input and frequency response

$$X(j\omega) \longrightarrow G(j\omega) \longrightarrow Y(j\omega) = G(j\omega)X(j\omega)$$

Fourier transform

Fourier transform (FT) and inverse FT

$$X(j\omega) = \int_{-\infty}^{\infty} x(t)e^{-j\omega t}dt \qquad x(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} X(j\omega)e^{j\omega t}d\omega$$

• FT of output is multiplication of FT of input and frequency response

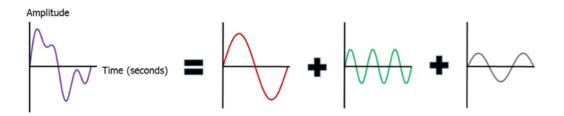
$$X(j\omega) \longrightarrow G(j\omega) \longrightarrow Y(j\omega) = G(j\omega)X(j\omega)$$

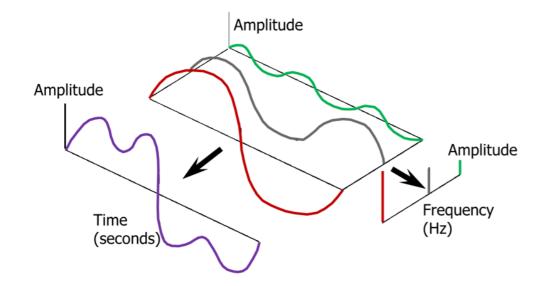
System output via inverse FT

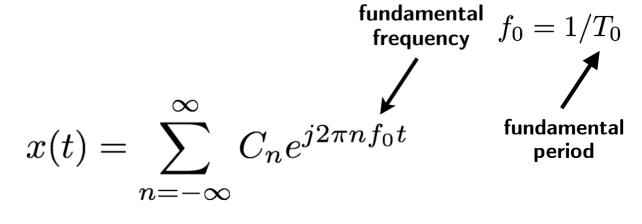
$$y(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} Y(j\omega) e^{j\omega t} d\omega$$

Fourier series

Fourier series (FS) for periodic signal



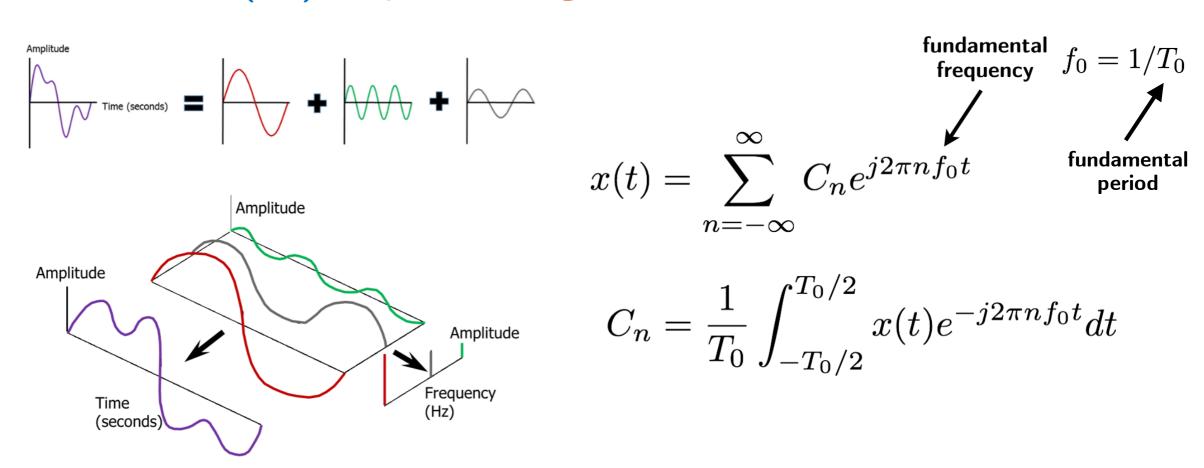




$$C_n = \frac{1}{T_0} \int_{-T_0/2}^{T_0/2} x(t)e^{-j2\pi n f_0 t} dt$$

Fourier series

Fourier series (FS) for periodic signal



 When the period approaches infinity, the spectrum becomes continuous leading to FT for aperiodic signal (previous slide)

Time domain vs Frequency domain

Theorem

Convolution in time domain is equivalent to multiplication in frequency domain, i.e.,

$$y(t) = g(t) * x(t) \equiv \mathcal{F}^{-1} \{ Y(j\omega) = G(j\omega)X(j\omega) \}$$

• Proof
$$\mathcal{F}\{g(t)*x(t)\} = \int_t \int_\tau g(t-\tau)x(\tau)d\tau e^{-j\omega t}dt$$

$$= \int_\tau x(\tau)e^{-j\omega\tau}d\tau \mathcal{F}\{g(t)\}$$

$$= \mathcal{F}\{g(t)\}\mathcal{F}\{x(t)\}$$

Time domain vs Frequency domain

Theorem

Convolution in time domain is equivalent to multiplication in frequency domain, i.e.,

$$y(t) = g(t) * x(t) \equiv \mathcal{F}^{-1} \{ Y(j\omega) = G(j\omega)X(j\omega) \}$$

• Proof
$$\mathcal{F}\{g(t)*x(t)\} = \int_t \int_\tau g(t-\tau)x(\tau)d\tau e^{-j\omega t}dt$$

$$= \int_\tau x(\tau)e^{-j\omega\tau}d\tau \mathcal{F}\{g(t)\}$$

$$= \mathcal{F}\{g(t)\}\mathcal{F}\{x(t)\}$$

Time domain vs Frequency domain

Theorem

Convolution in time domain is equivalent to multiplication in frequency domain, i.e.,

$$y(t) = g(t) * x(t) \equiv \mathcal{F}^{-1}\{Y(j\omega) = G(j\omega)X(j\omega)\}\$$

• Proof
$$\mathcal{F}\{g(t)*x(t)\} = \int_{t}^{t} \int_{\tau}^{t} g(t-\tau)x(\tau)d\tau e^{-j\omega t}dt$$

$$= \int_{\tau}^{t} x(\tau)e^{-j\omega \tau}d\tau \mathcal{F}\{g(t)\}$$

$$= \mathcal{F}\{g(t)\}\mathcal{F}\{x(t)\}$$

Time domain vs Frequency domain

Theorem

Convolution in time domain is equivalent to multiplication in frequency domain, i.e.,

$$y(t) = g(t) * x(t) \equiv \mathcal{F}^{-1} \{ Y(j\omega) = G(j\omega)X(j\omega) \}$$

$$\begin{aligned} \bullet \ \, \mathbf{Proof} \quad & \mathcal{F}\{g(t)*x(t)\} = \int_t \int_\tau g(t-\tau)x(\tau)d\tau e^{-j\omega t}dt \\ = & \int_\tau x(\tau)e^{-j\omega\tau}d\tau \mathcal{F}\{g(t)\} \\ & = & \mathcal{F}\{g(t)\} \boxed{\mathcal{F}\{x(t)\}} \end{aligned}$$

Time domain vs Frequency domain

Theorem

Convolution in time domain is equivalent to multiplication in frequency domain, i.e.,

$$y(t) = g(t) * x(t) \equiv \mathcal{F}^{-1}\{Y(j\omega) = G(j\omega)X(j\omega)\}\$$

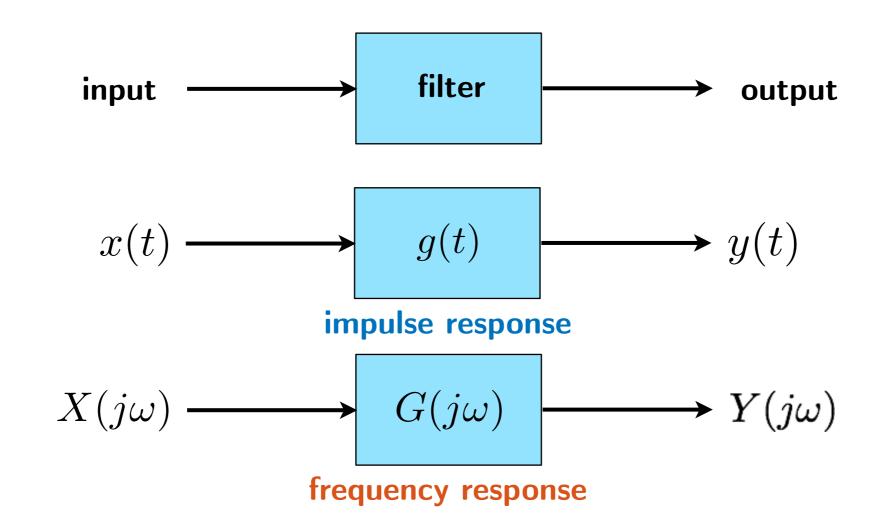
• Proof
$$\mathcal{F}\{g(t)*x(t)\} = \int_t \int_\tau g(t-\tau)x(\tau)d\tau e^{-j\omega t}dt$$

$$= \int_\tau x(\tau)e^{-j\omega\tau}d\tau \mathcal{F}\{g(t)\}$$

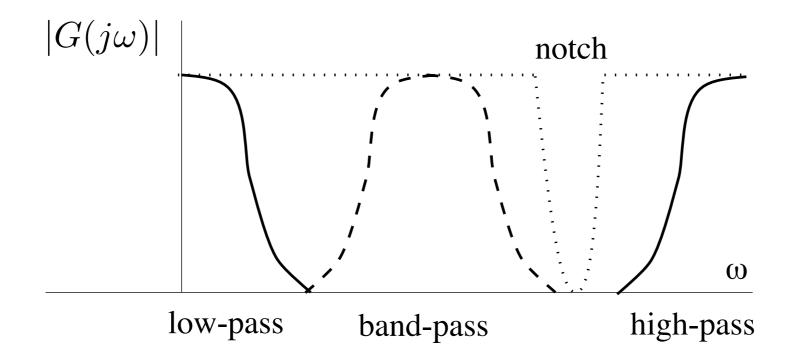
$$= \mathcal{F}\{g(t)\}\mathcal{F}\{x(t)\}$$

This allows us to move losslessly between **time** and **frequency** domains, choosing whichever is the easier to work with

Filtering as input-output relationship



- Filters are **frequency-selective** linear systems
 - low-pass: extract average or eliminate high-frequency fluctuations
 - high-pass: follow small-amplitude high-frequency perturbations in presence of much larger slowly-varying component

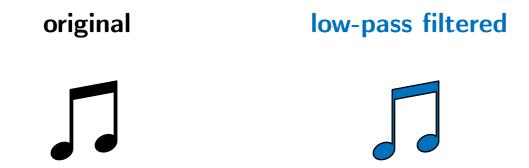


- Filters are **frequency-selective** linear systems
 - low-pass: extract average or eliminate high-frequency fluctuations
 - high-pass: follow small-amplitude high-frequency perturbations in presence of much larger slowly-varying component

original



- Filters are **frequency-selective** linear systems
 - low-pass: extract average or eliminate high-frequency fluctuations
 - high-pass: follow small-amplitude high-frequency perturbations in presence of much larger slowly-varying component



- Filters are **frequency-selective** linear systems
 - low-pass: extract average or eliminate high-frequency fluctuations
 - high-pass: follow small-amplitude high-frequency perturbations in presence of much larger slowly-varying component



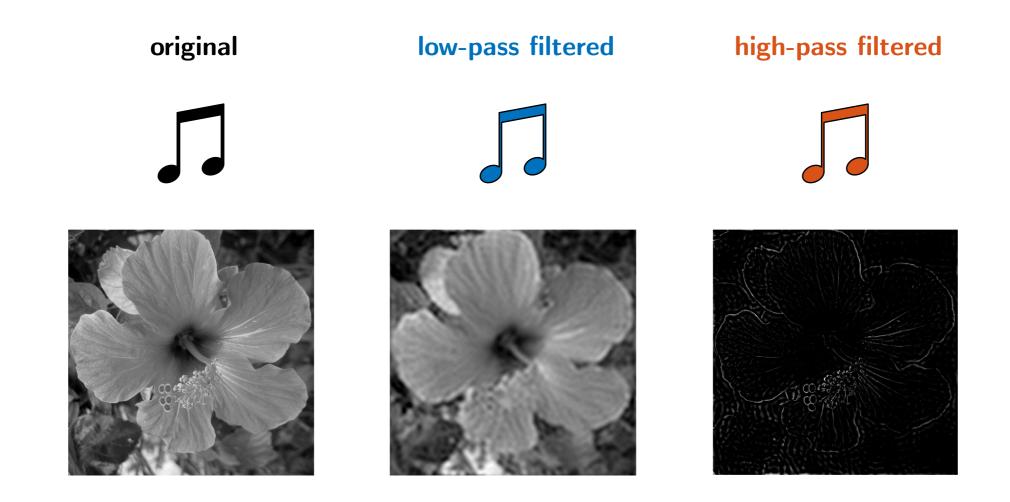
- Filters are frequency-selective linear systems
 - low-pass: extract average or eliminate high-frequency fluctuations
 - high-pass: follow small-amplitude high-frequency perturbations in presence of much larger slowly-varying component



- Filters are **frequency-selective** linear systems
 - low-pass: extract average or eliminate high-frequency fluctuations
 - high-pass: follow small-amplitude high-frequency perturbations in presence of much larger slowly-varying component



- Filters are frequency-selective linear systems
 - low-pass: extract average or eliminate high-frequency fluctuations
 - high-pass: follow small-amplitude high-frequency perturbations in presence of much larger slowly-varying component



Fourier transform:
$$\hat{f}(\omega) = \int (e^{j\omega x})^* f(x) dx$$
 $f(x) = \frac{1}{2\pi} \int \hat{f}(\omega) e^{j\omega x} d\omega$

Fourier transform:
$$\hat{f}(\omega) = \int (e^{j\omega x})^* f(x) dx$$
 $f(x) = \frac{1}{2\pi} \int \hat{f}(\omega) e^{j\omega x} d\omega$

f

Fourier transform:
$$\hat{f}(\omega) = \int (e^{j\omega x})^* f(x) dx$$
 $f(x) = \frac{1}{2\pi} \int \hat{f}(\omega) e^{j\omega x} d\omega$

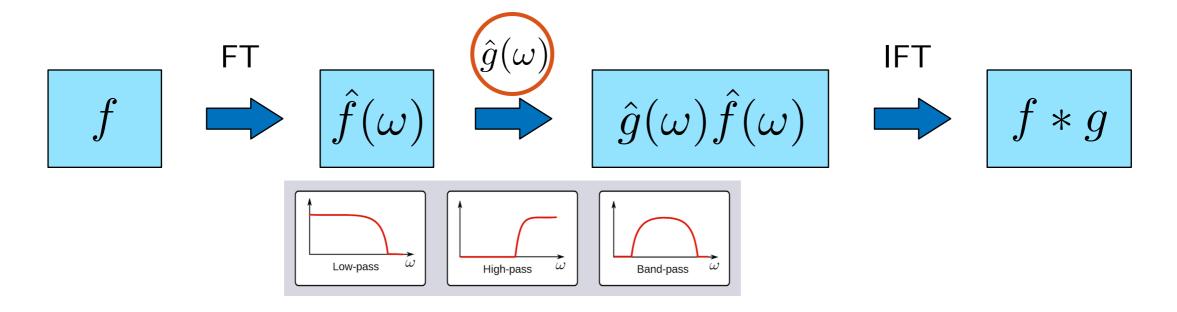
$$f$$
 $\hat{f}(\omega)$

Fourier transform:
$$\hat{f}(\omega) = \int (e^{j\omega x})^* f(x) dx$$
 $f(x) = \frac{1}{2\pi} \int \hat{f}(\omega) e^{j\omega x} d\omega$

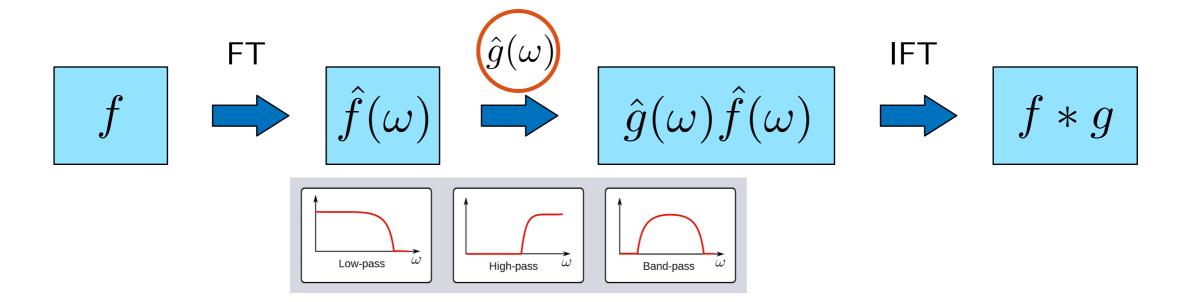
$$\begin{array}{c|c} & & \hat{g}(\omega) \\ \hline f & & & \\ \hline \end{array} \qquad \begin{array}{c|c} \hat{g}(\omega) \\ \hline & & \\ \hline & & \\ \hline \end{array} \qquad \begin{array}{c|c} \hat{g}(\omega) \\ \hline & & \\ \hline & \hat{g}(\omega) \hat{f}(\omega) \\ \hline \end{array}$$

Fourier transform:
$$\hat{f}(\omega) = \int (e^{j\omega x})^* f(x) dx$$
 $f(x) = \frac{1}{2\pi} \int \hat{f}(\omega) e^{j\omega x} d\omega$

Fourier transform:
$$\hat{f}(\omega) = \int (e^{j\omega x})^* f(x) dx$$
 $f(x) = \frac{1}{2\pi} \int \hat{f}(\omega) e^{j\omega x} d\omega$



Fourier transform:
$$\hat{f}(\omega) = \int (e^{j\omega x})^* f(x) dx$$
 $f(x) = \frac{1}{2\pi} \int \hat{f}(\omega) e^{j\omega x} d\omega$



Takeaway: signal processing (filtering) requires two considerations

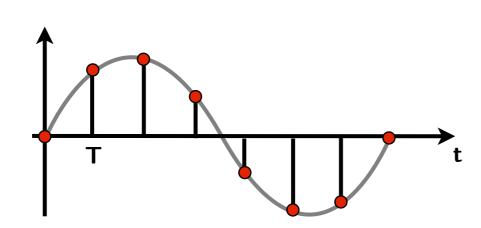
- find a **good representation** (e.g., via Fourier transform) of the signal
- design (or learn) appropriate filters (note that filtering = convolution!)

Lecture 1

- Introduction & Basic concepts and tools
- A historical overview of signal representation techniques
- Applications & Discussion

What is good representation of a signal?

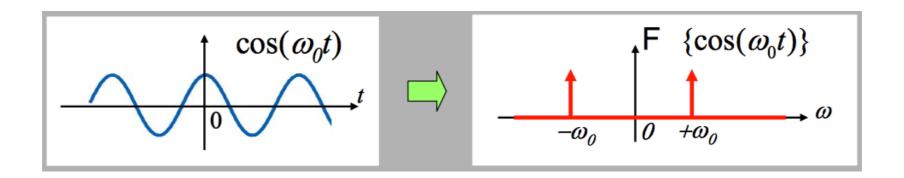
- Sum of delta functions in time or space (sampling domain)
 - good for display or playback
 - not good for analysis (e.g., denoising, compression)





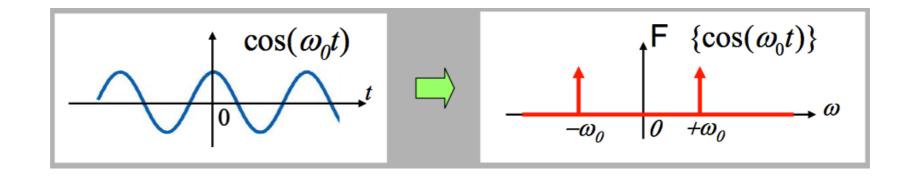
What is good representation of a signal?

- Representation often involves transformation of the signal into a new domain where signal characteristics are revealed
 - example: Fourier coefficients reveal rate of change of the signal



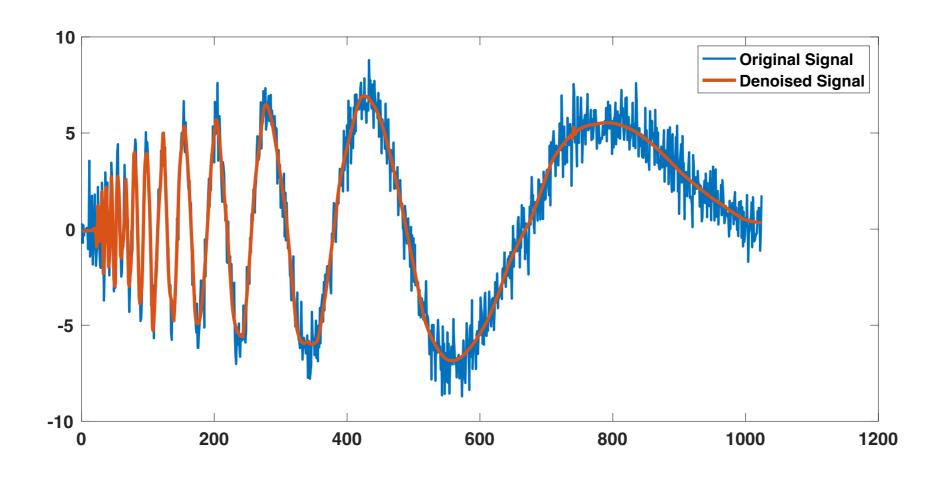
What is good representation of a signal?

- Representation often involves transformation of the signal into a new domain where signal characteristics are revealed
 - example: Fourier coefficients reveal rate of change of the signal



- Usefulness of the representation depends on the analysis goal
 - which may vary but all shares the core desire for simplification

Example: Denoising



goal: recover signal from noisy observation

Example: Compression

original



JPEG 2000 (10% in size)

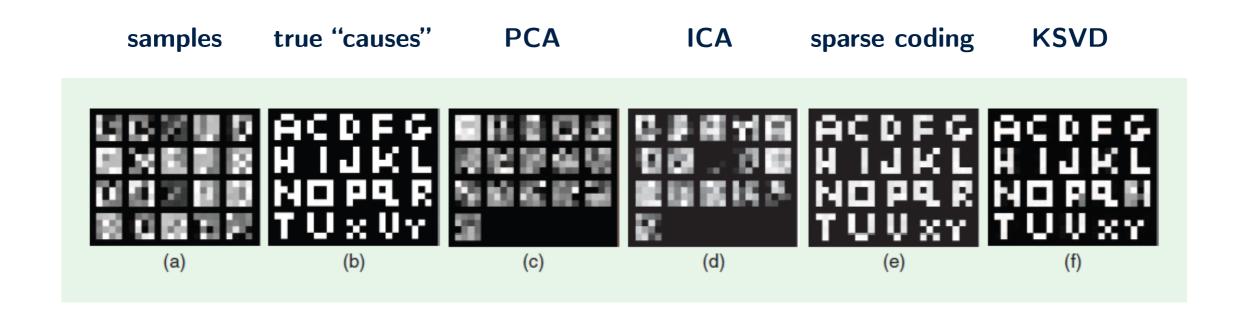


JPEG 2000 (1% in size)

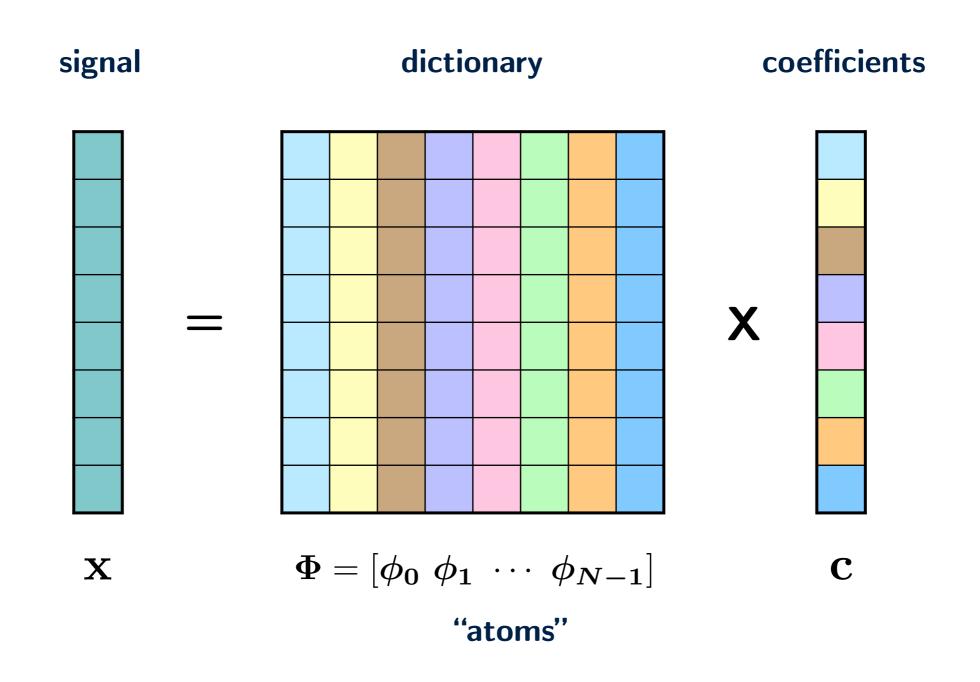


goal: compress signal without sacrificing quality

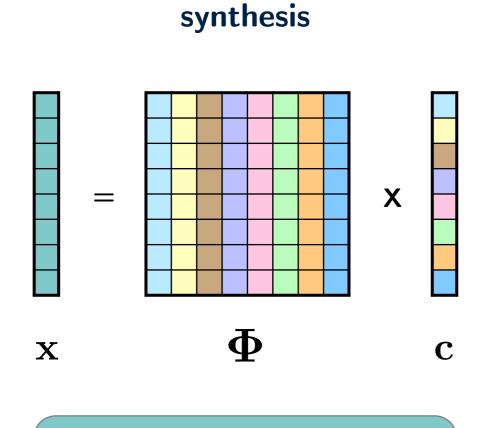
Example: Recognition



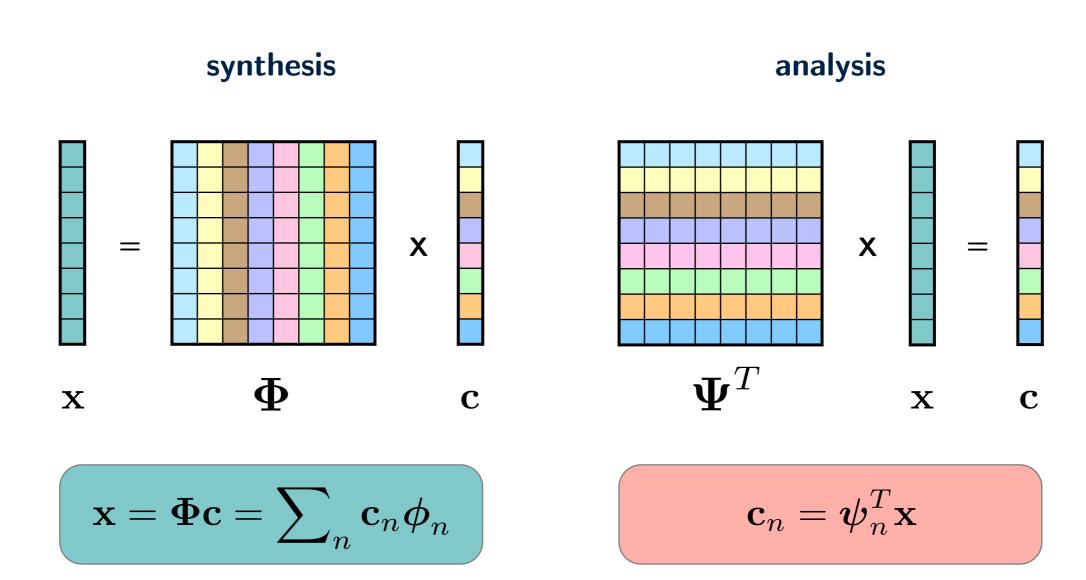
goal: capture true "causes" of signal



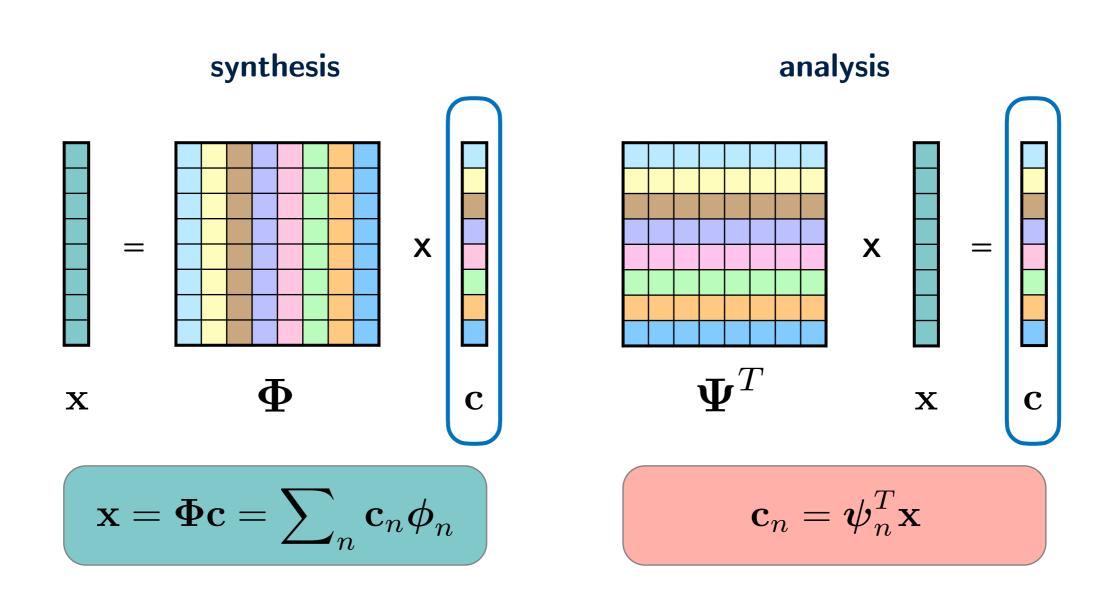
Complete dictionaries



Complete dictionaries

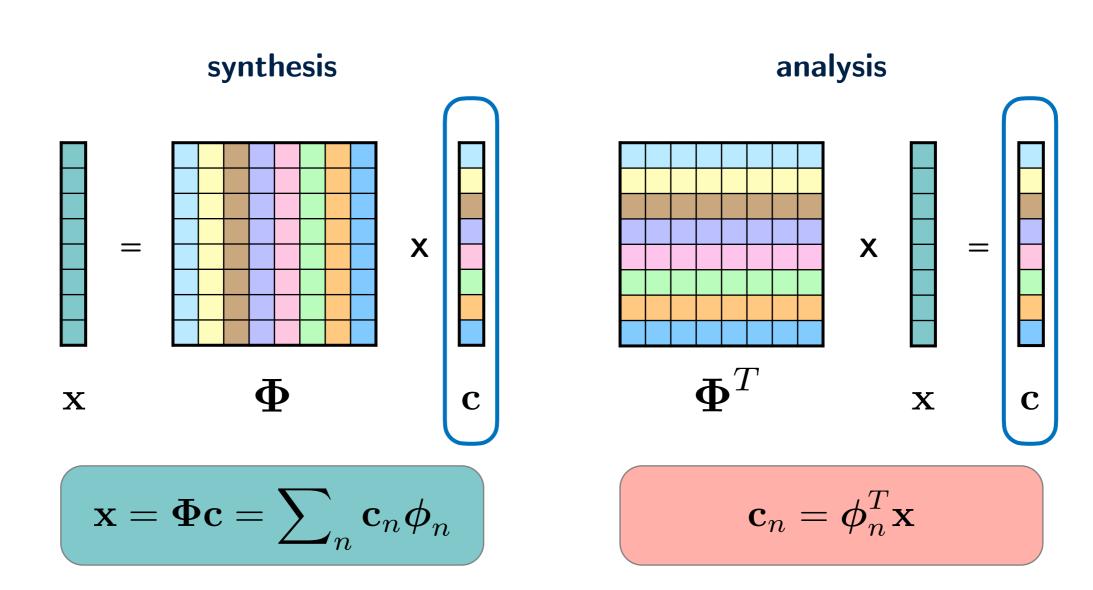


Complete dictionaries



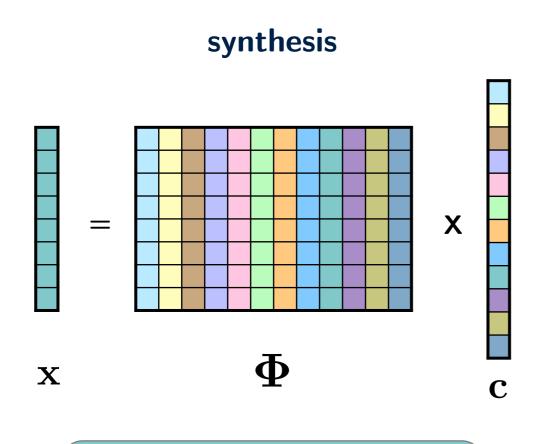
equivalent for complete and biorthogonal dictionaries

Complete dictionaries



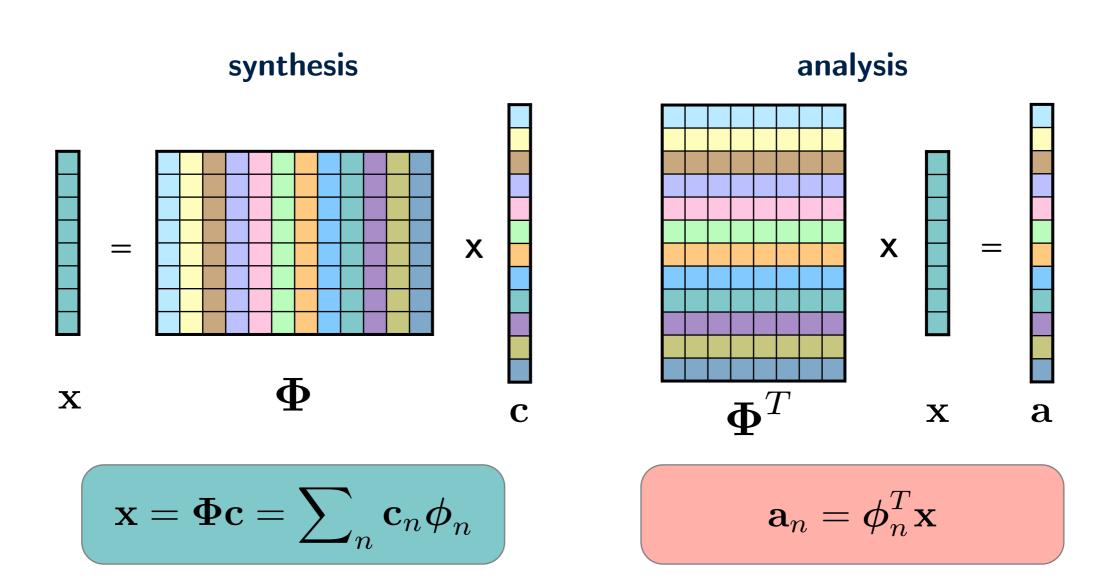
equivalent for complete and biorthogonal dictionaries

• Overcomplete (redundant) dictionaries

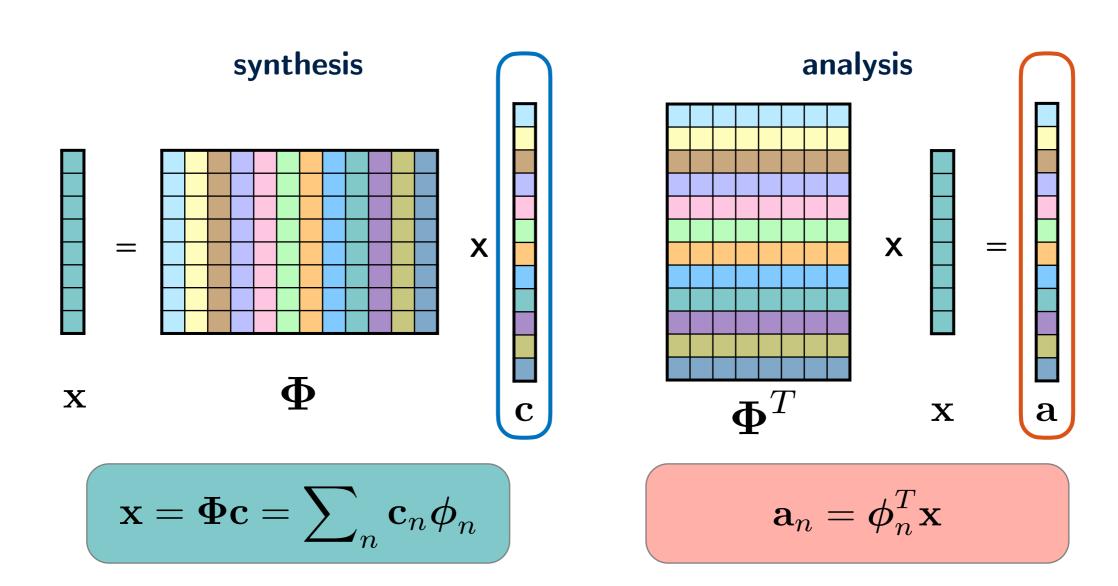


$$\mathbf{x} = \mathbf{\Phi} \mathbf{c} = \sum_n \mathbf{c}_n \boldsymbol{\phi}_n$$

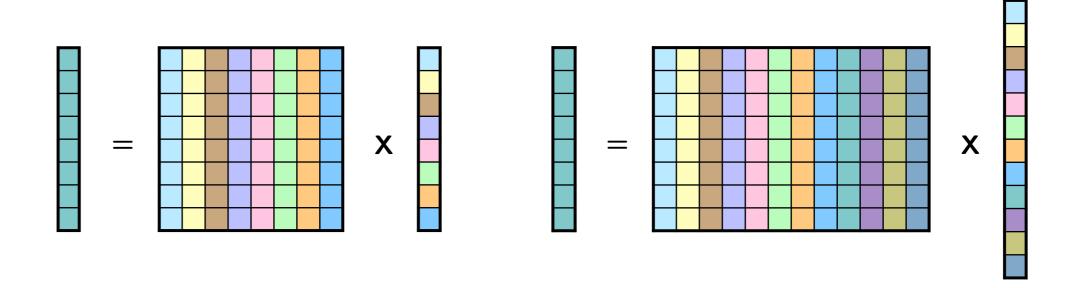
• Overcomplete (redundant) dictionaries



Overcomplete (redundant) dictionaries



not equivalent for overcomplete dictionaries



Two sources of dictionary design

- mathematical modelling of data (transforms/analytic dictionaries)
- a set of realisations of data (learned dictionaries)

1960s: Fourier basis and DFT



- recall the LTI system

$$x(t) \longrightarrow g(t) \longrightarrow y(t)$$

1960s: Fourier basis and DFT



recall the LTI system

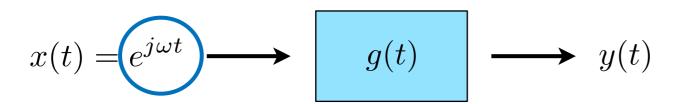
$$x(t) = e^{j\omega t} \longrightarrow g(t) \longrightarrow y(t)$$

$$y(t) = \int_{-\infty}^{\infty} e^{j\omega(t-\tau)} g(\tau) d\tau = e^{j\omega t} G(\omega)$$

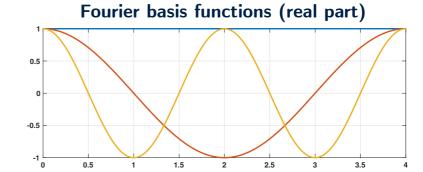
1960s: Fourier basis and DFT



recall the LTI system



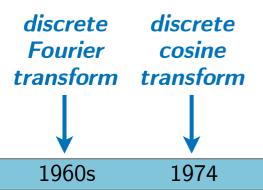
 $y(t) = \int_{-\infty}^{\infty} e^{j\omega(t-\tau)} g(\tau) d\tau = e^{j\omega t} G(\omega)$



Fourier basis diagonalises convolution operator

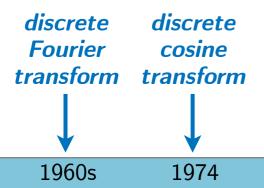
$$X(\omega) \longrightarrow G(\omega) \longrightarrow Y(\omega) = X(\omega)G(\omega)$$

1960s: Fourier basis and DFT



- Fourier basis describes a signal in terms of its **global** frequency content and hence is good at representing **uniformly smooth** signals

1960s: Fourier basis and DFT

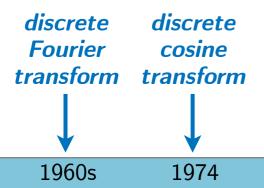


- Fourier basis describes a signal in terms of its global frequency content and hence is good at representing uniformly smooth signals
- discrete Fourier transform (DFT) provides an orthogonal dictionary: $\phi_n(k) = e^{j \frac{2\pi}{N} nk}$

$$\begin{pmatrix} x[0] \\ x[1] \\ x[2] \\ \vdots \\ x[N-1] \end{pmatrix} = \frac{1}{N} \begin{pmatrix} 1 & 1 & 1 & 1 & \dots & 1 \\ 1 & W & W^2 & W^3 & \dots & W^{N-1} \\ 1 & W^2 & W^4 & W^6 & \dots & W^{N-2} \\ 1 & W^3 & W^6 & W^9 & \dots & W^{N-3} \\ \vdots & & & & & & \\ 1 & W^{N-1} & W^{N-2} & W^{N-3} & \dots & W \end{pmatrix} \begin{pmatrix} X[0] \\ X[1] \\ X[2] \\ \vdots \\ X[N-1] \end{pmatrix} \quad \text{with} \quad W = e^{j\frac{2\pi}{N}}$$

- fast Fourier transform (FFT) reduces complexity from $\mathcal{O}(N^2)$ to $\mathcal{O}(N{\log}N)$

1960s: Fourier basis and DFT

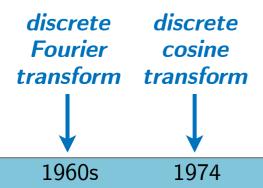


- Fourier basis describes a signal in terms of its global frequency content and hence is good at representing uniformly smooth signals
- discrete Fourier transform (DFT) provides an orthogonal dictionary: $\phi_n(k) = e^{j \frac{2\pi}{N} nk}$

$$\begin{pmatrix} x[0] \\ x[1] \\ x[2] \\ \vdots \\ x[N-1] \end{pmatrix} = \frac{1}{N} \begin{pmatrix} 1 & 1 & 1 & 1 & \dots & 1 \\ 1 & W & W^2 & W^3 & \dots & W^{N-1} \\ 1 & W^2 & W^4 & W^6 & \dots & W^{N-2} \\ 1 & W^3 & W^6 & W^9 & \dots & W^{N-3} \\ \vdots & & & & & & \\ 1 & W^{N-1} & W^{N-2} & W^{N-3} & \dots & W \end{pmatrix} \begin{pmatrix} X[0] \\ X[1] \\ X[2] \\ \vdots \\ X[N-1] \end{pmatrix} \quad \text{with} \quad W = e^{j\frac{2\pi}{N}}$$

- fast Fourier transform (FFT) reduces complexity from $\mathcal{O}(N^2)$ to $\mathcal{O}(N \log N)$
- can be made into a real transform called discrete cosine transform (DCT)

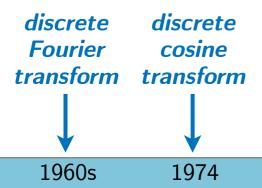
The need for sparsity



projection onto a fixed subset of DFT/DCT atoms leads to compaction

$$\mathbf{x} pprox \sum_{n \in \mathcal{S}_k} (\mathbf{\Psi}_n^T \mathbf{x}) \mathbf{\Phi}_n$$

The need for sparsity



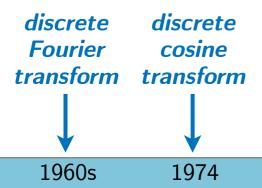
projection onto a fixed subset of DFT/DCT atoms leads to compaction

$$\mathbf{x} pprox \sum_{n \in \mathcal{S}_k} (\mathbf{\Psi}_n^T \mathbf{x}) \mathbf{\Phi}_n$$

- from compaction (simplicity) to sparsity: signal as linear combination of a few atoms
- sparsity requires shift from linear to nonlinear approximation

$$\mathbf{x} \approx \sum_{k} c_k \phi_k$$
 subset of atoms (different for each x)

The need for sparsity



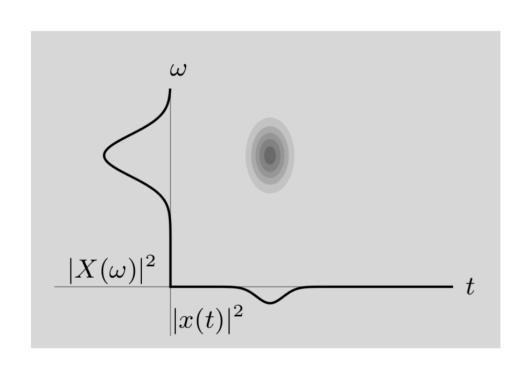
projection onto a fixed subset of DFT/DCT atoms leads to compaction

$$\mathbf{x} pprox \sum_{n \in \mathcal{S}_k} (\mathbf{\Psi}_n^T \mathbf{x}) \mathbf{\Phi}_n$$

- from compaction (simplicity) to sparsity: signal as linear combination of a few atoms
- sparsity requires shift from linear to nonlinear approximation

$$\mathbf{x} \approx \sum_{k} c_k \phi_k$$
 subset of atoms (different for each x)

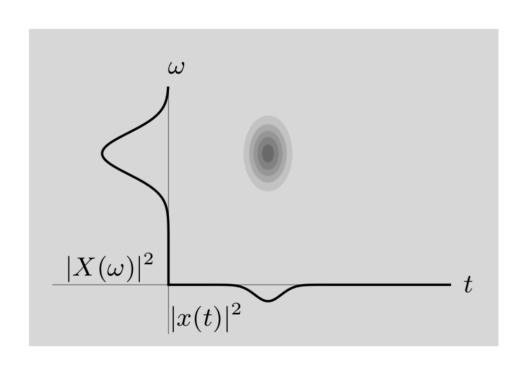
- sparsity requires localisation: atoms with concentrated support
 - allow more flexible representations based on local characteristics
 - limit effects of irregularities (a main source of large coefficients)



time localisation

$$\mu_t = \frac{1}{||x||^2} \int_{-\infty}^{\infty} t |x(t)|^2 dt$$

$$\Delta_t = \left(\frac{1}{||x||^2} \int_{-\infty}^{\infty} (t - \mu_t)^2 |x(t)|^2 dt\right)^{\frac{1}{2}}$$



time localisation

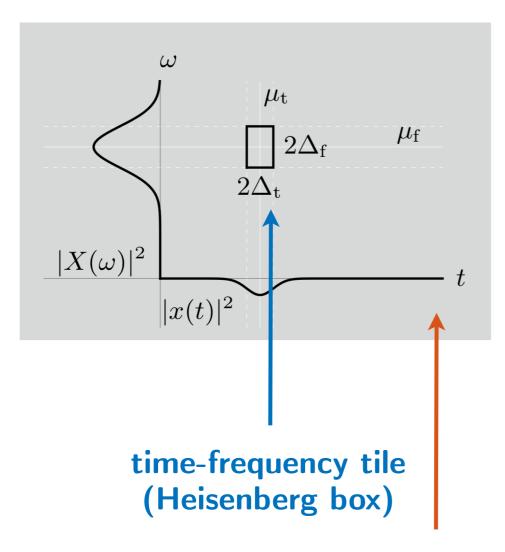
$$\mu_t = \frac{1}{||x||^2} \int_{-\infty}^{\infty} t |x(t)|^2 dt$$

$$\Delta_t = \left(\frac{1}{||x||^2} \int_{-\infty}^{\infty} (t - \mu_t)^2 |x(t)|^2 dt\right)^{\frac{1}{2}}$$

frequency localisation

$$\mu_f = \frac{1}{2\pi||x||^2} \int_{-\infty}^{\infty} \omega X(\omega)|^2 d\omega$$

$$\Delta_f = \left(\frac{1}{2\pi||x||^2} \int_{-\infty}^{\infty} (\omega - \mu_f)^2 |X(\omega)|^2 d\omega\right)^{\frac{1}{2}}$$



time-frequency plane

time localisation

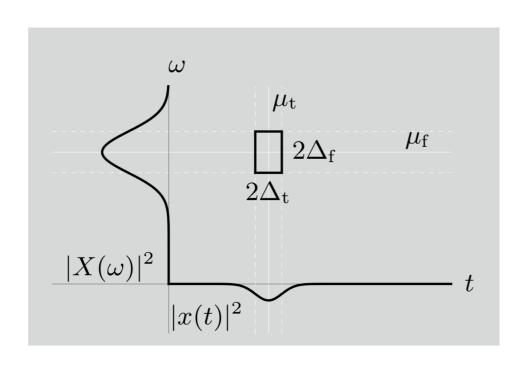
$$\mu_t = \frac{1}{\|x\|^2} \int_{-\infty}^{\infty} t |x(t)|^2 dt$$

$$\Delta_t = \left(\frac{1}{||x||^2} \int_{-\infty}^{\infty} (t - \mu_t)^2 |x(t)|^2 dt\right)^{\frac{1}{2}}$$

frequency localisation

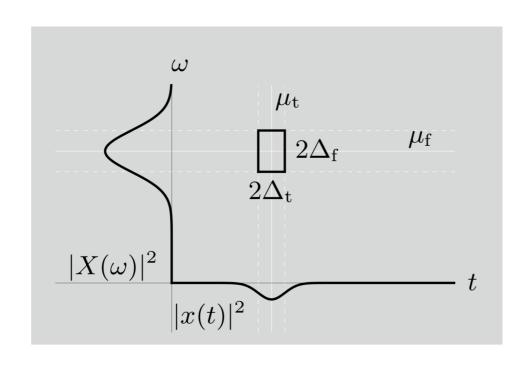
$$\mu_f = \frac{1}{2\pi||x||^2} \int_{-\infty}^{\infty} \omega X(\omega)|^2 d\omega$$

$$\Delta_f = \left(\frac{1}{2\pi||x||^2} \int_{-\infty}^{\infty} (\omega - \mu_f)^2 |X(\omega)|^2 d\omega\right)^{\frac{1}{2}}$$



Heisenberg's uncertainty principle

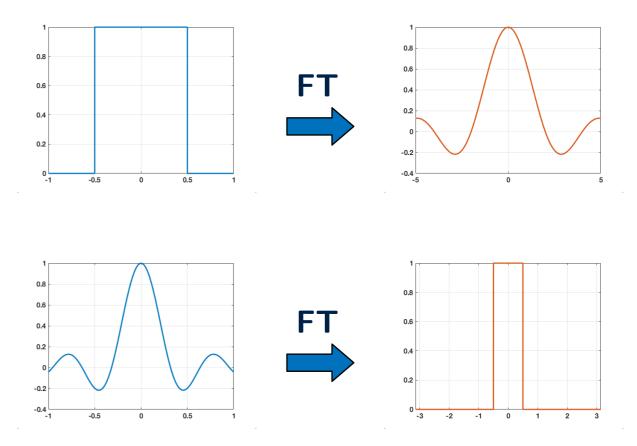
Let
$$x \in \mathcal{L}^2(\mathbb{R})$$
, then $\Delta_t \Delta_f \ge \frac{1}{2}$



Heisenberg's uncertainty principle

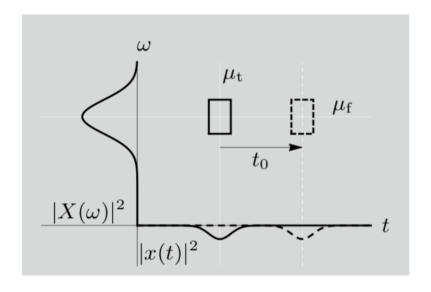
Let
$$x \in \mathcal{L}^2(\mathbb{R})$$
, then $\Delta_t \Delta_f \ge \frac{1}{2}$

examples



Consider three basic operations

shift in time



$$y(t) = x(t - t_0)$$

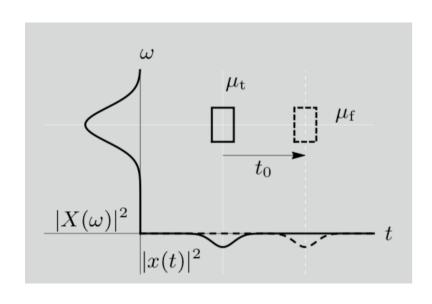


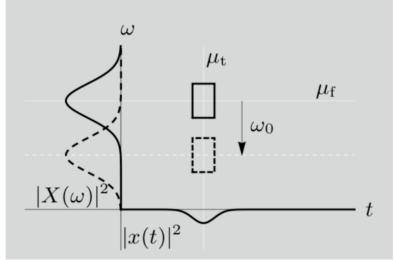
$$Y(\omega) = e^{-j\omega t_0} X(\omega)$$

Consider three basic operations

shift in time

shift in frequency





$$y(t) = x(t - t_0)$$



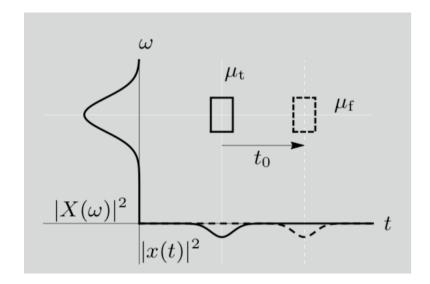
$$Y(\omega) = e^{-j\omega t_0} X(\omega)$$

$$y(t) = e^{j\omega_0 t} x(t)$$

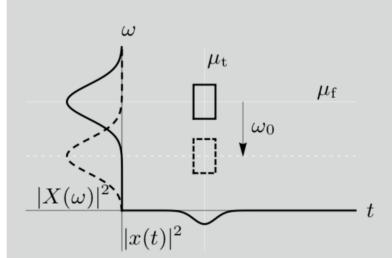
$$Y(\omega) = X(\omega - \omega_0)$$

Consider three basic operations

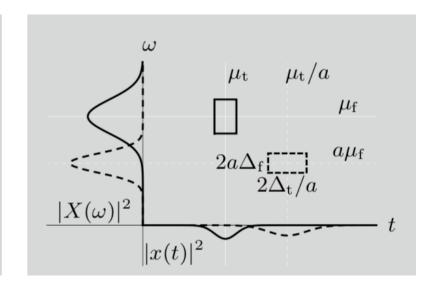
shift in time



shift in frequency



scaling in time



$$y(t) = x(t - t_0)$$



$$Y(\omega) = e^{-j\omega t_0} X(\omega)$$

$$y(t) = e^{j\omega_0 t} x(t)$$



$$Y(\omega) = X(\omega - \omega_0)$$

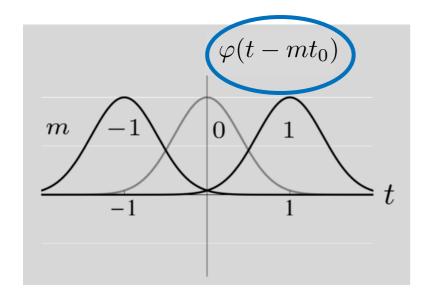
$$y(t) = \sqrt{a}x(at)$$

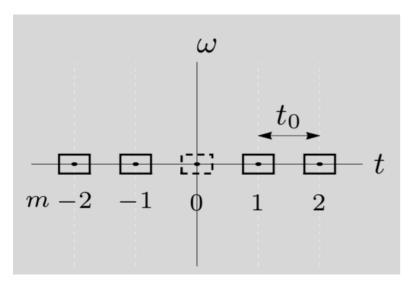


$$Y(\omega) = \frac{1}{\sqrt{a}}X(\frac{\omega}{a})$$

Consider three structured sets of functions

shift in time

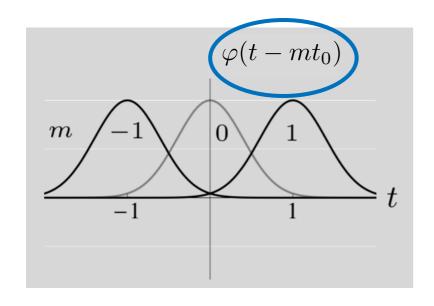


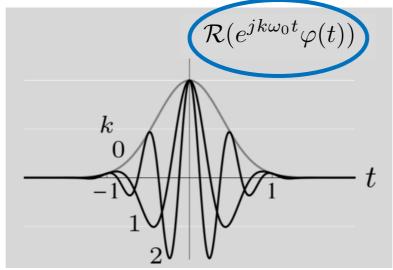


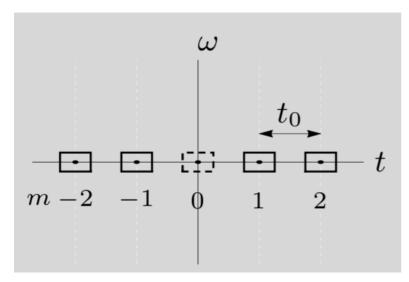
Consider three structured sets of functions

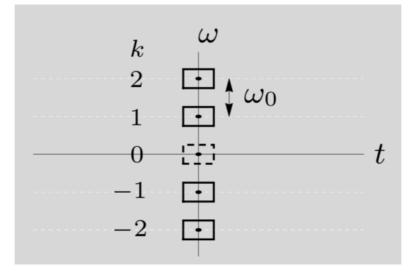
shift in time

shift in frequency



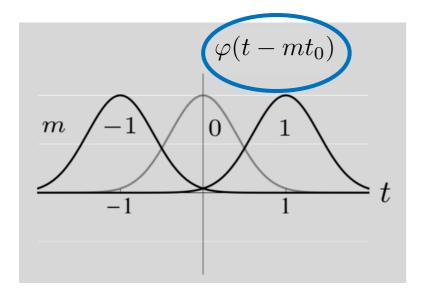




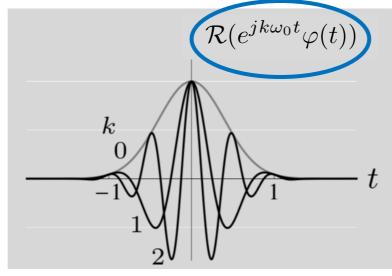


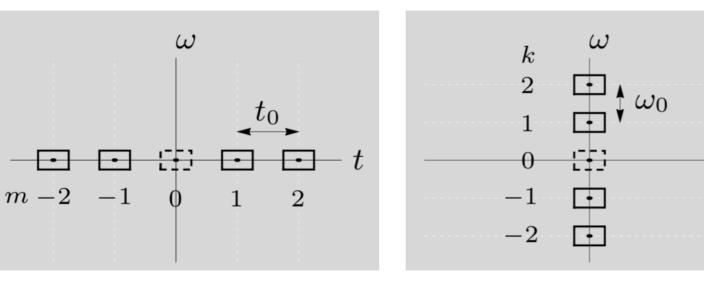
Consider three structured sets of functions

shift in time

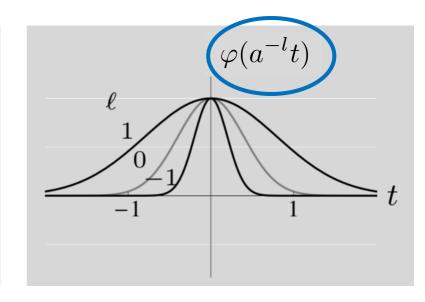


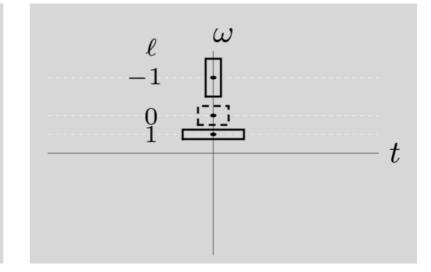
shift in frequency



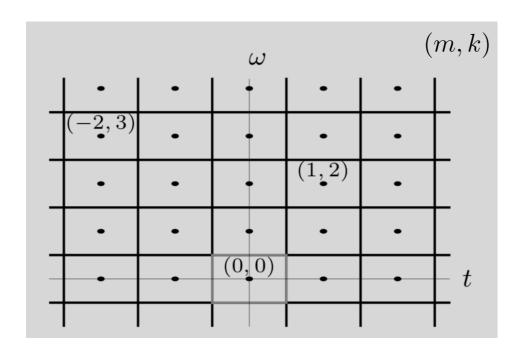


scaling in time



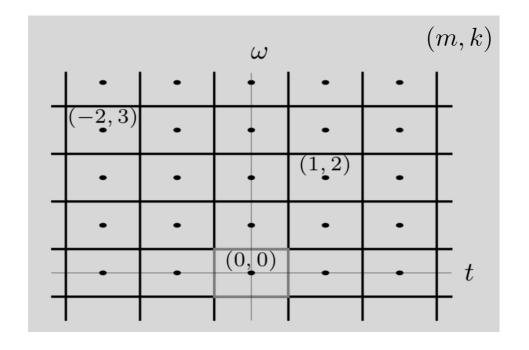


time shift and modulation

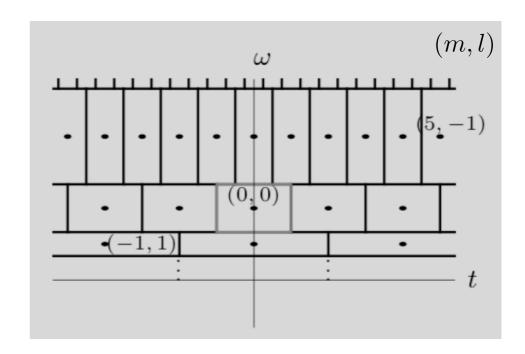


$$\varphi_{k,m}(t) = e^{jk\omega_0 t} \varphi(t - mt_0) \quad k, m \in \mathbb{Z}$$

time shift and modulation

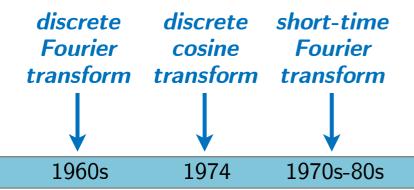


time shift and scaling

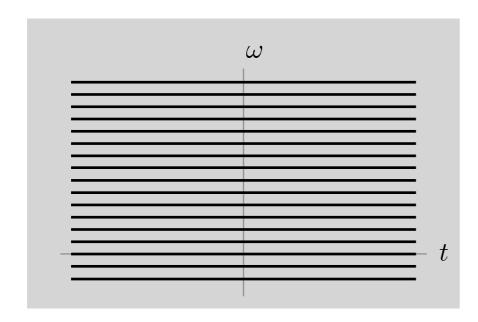


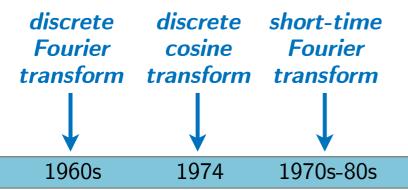
$$\varphi_{k,m}(t) = e^{jk\omega_0 t} \varphi(t - mt_0) \quad k, m \in \mathbb{Z} \qquad \varphi_{l,m}(t) = \varphi(a^{-l}t - mt_0) \quad l, m \in \mathbb{Z}$$

$$\varphi_{l,m}(t) = \varphi(a^{-l}t - mt_0) \quad l, m \in \mathbb{Z}$$
$$= \varphi(a^{-l}(t - ma^l t_0))$$



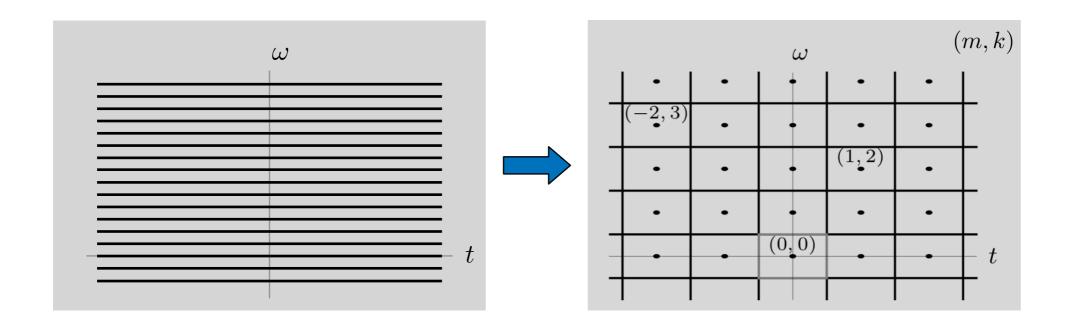
- Fourier basis functions (complex exponentials) are not localised in time

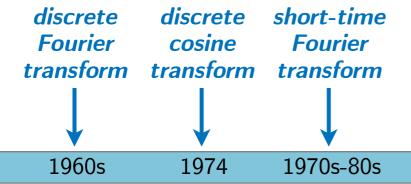




- Fourier basis functions (complex exponentials) are not localised in time
- consider a set of shifted and modulated versions of a low-pass function

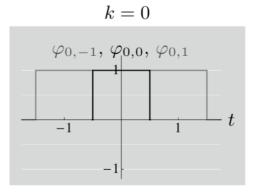
$$\varphi_{k,m}(t) = e^{jk\omega_0 t} \varphi(t - mt_0) \quad k, m \in \mathbb{Z}$$

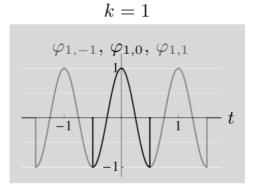


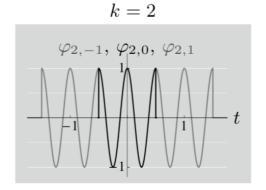


- example: consider a box function and $t_0=1$, $\omega_0=2\pi$

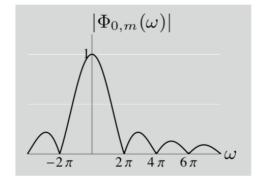
$$\varphi_{k,m}(t) = e^{jk2\pi t} \varphi(t-m), \quad \varphi(t) = \begin{cases} 1, & \text{for } |t| \leq \frac{1}{2}; \\ 0, & \text{otherwise.} \end{cases}$$

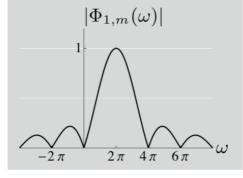


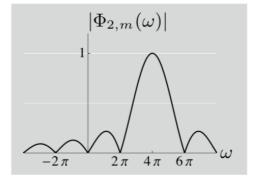




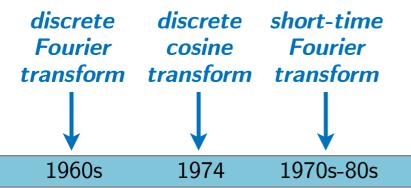
Basis functions (real parts only).





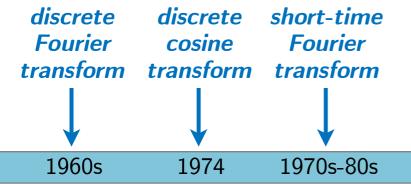


Magnitudes of the Fourier transform.



- we can define the following transform

$$X_{k,m} = \int_{-\infty}^{\infty} x(t)\varphi_{k,m}^*(t)dt = \int_{-\infty}^{\infty} x(t)\varphi(t - mt_0)e^{-jk\omega_0 t}dt$$

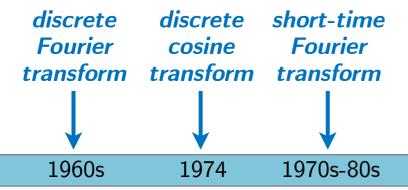


we can define the following transform

$$X_{k,m} = \int_{-\infty}^{\infty} x(t)\varphi_{k,m}^*(t)dt = \int_{-\infty}^{\infty} x(t)\varphi(t - mt_0)e^{-jk\omega_0 t}dt$$



$$X(\omega,\tau) = \int_{-\infty}^{\infty} x(t)\varphi_{\omega,\tau}^*(t)dt = \int_{-\infty}^{\infty} x(t)\varphi(t-\tau)e^{-j\omega t}dt$$

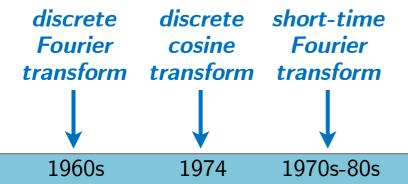


- we can define the following transform

$$X_{k,m} = \int_{-\infty}^{\infty} x(t)\varphi_{k,m}^*(t)dt = \int_{-\infty}^{\infty} x(t)\varphi(t - mt_0)e^{-jk\omega_0 t}dt$$

$$X(\omega,\tau) = \int_{-\infty}^{\infty} x(t)\varphi_{\omega,\tau}^{*}(t)dt = \int_{-\infty}^{\infty} x(t)\varphi(t-\tau)e^{-j\omega t}dt$$

applying time-localised window to the signal before taking Fourier transform:
 windowed or short-time Fourier transform (STFT)



- we can define the following transform

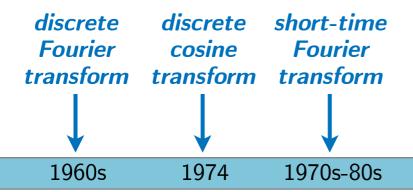
$$X_{k,m} = \int_{-\infty}^{\infty} x(t)\varphi_{k,m}^*(t)dt = \int_{-\infty}^{\infty} x(t)\varphi(t - mt_0)e^{-jk\omega_0 t}dt$$



$$X(\omega,\tau) = \int_{-\infty}^{\infty} x(t)\varphi_{\omega,\tau}^{*}(t)dt = \int_{-\infty}^{\infty} x(t)\varphi(t-\tau)e^{-j\omega t}dt$$

- applying time-localised window to the signal before taking Fourier transform:
 windowed or short-time Fourier transform (STFT)
- Gaussian window achieves localisation in frequency: Gabor transform
- STFT maps a 1-D function into a 2-D function (overcomplete)

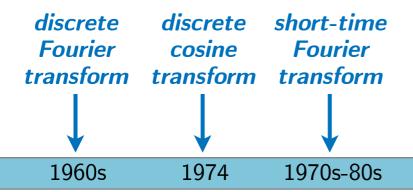
DCT vs STFT



- discrete STFT generally provides an overcomplete dictionary

$$\phi_{k,m}(n) = e^{j\frac{2\pi}{N}nk}\varphi(n - mN)$$

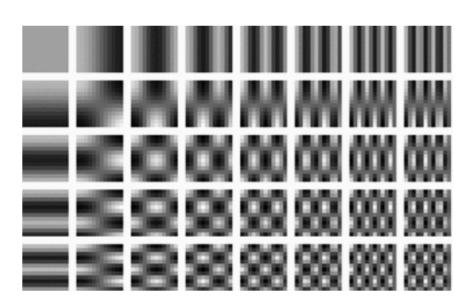
DCT vs STFT



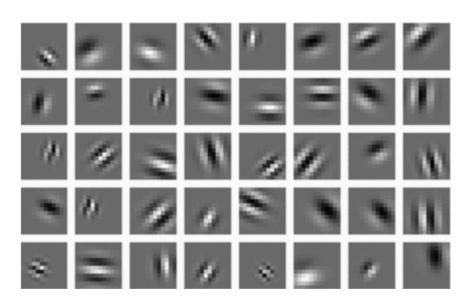
- discrete STFT generally provides an overcomplete dictionary

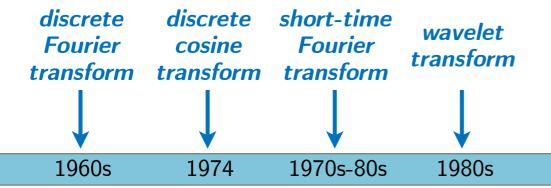
$$\phi_{k,m}(n) = e^{j\frac{2\pi}{N}nk}\varphi(n - mN)$$

DCT

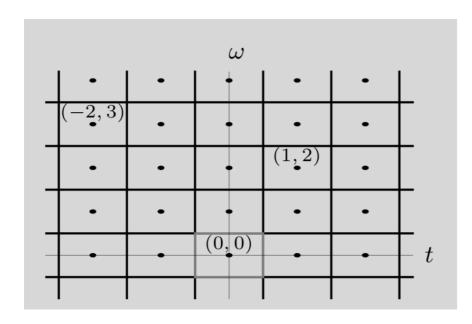


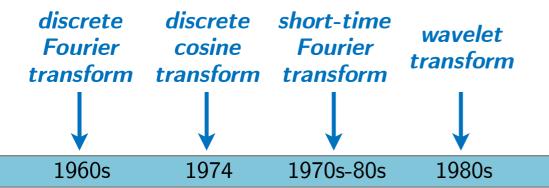
STFT



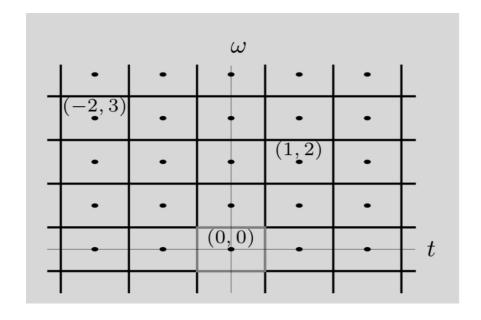


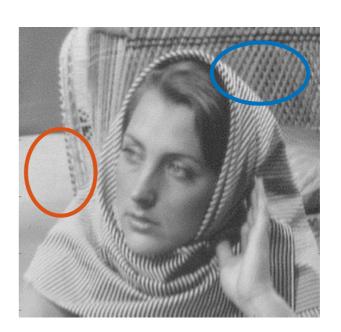
- STFT atoms have fixed time-frequency resolution

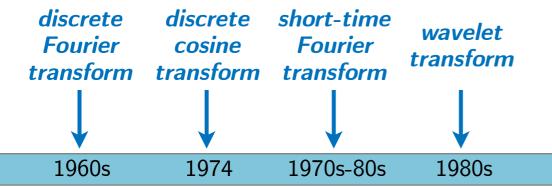




- STFT atoms have fixed time-frequency resolution
- often times a multiresolution representation is needed to capture various scales in natural signals

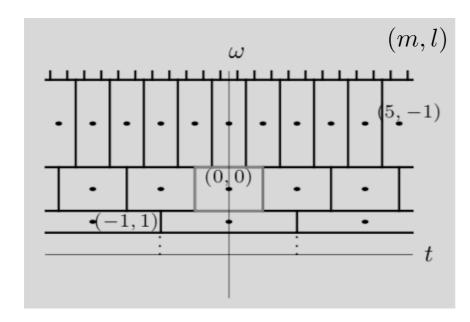


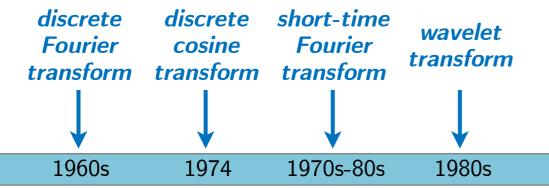




consider a set of shifted and scaled versions of a band-pass function

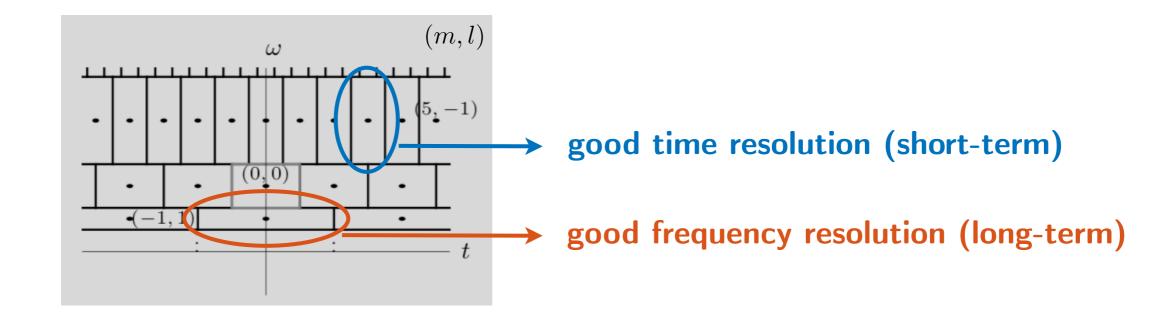
$$\varphi_{l,m}(t) = \varphi(a^{-l}t - mt_0) = \varphi(\frac{t - ma^l t_0}{a^l}) \quad l, m \in \mathbb{Z}$$

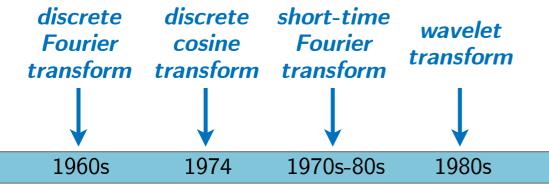




consider a set of shifted and scaled versions of a band-pass function

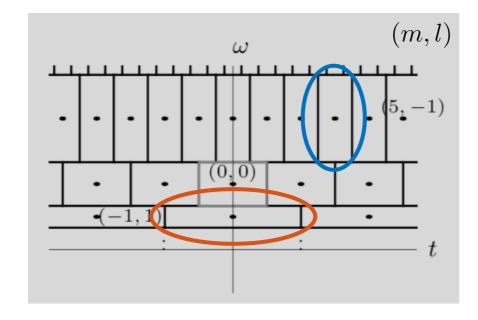
$$\varphi_{l,m}(t) = \varphi(a^{-l}t - mt_0) = \varphi(\frac{t - ma^l t_0}{a^l}) \quad l, m \in \mathbb{Z}$$

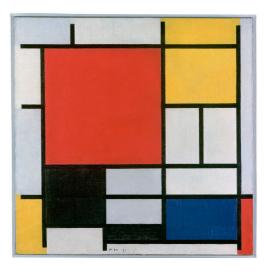




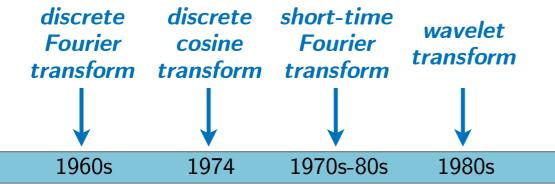
consider a set of shifted and scaled versions of a band-pass function

$$\varphi_{l,m}(t) = \varphi(a^{-l}t - mt_0) = \varphi(\frac{t - ma^l t_0}{a^l}) \quad l, m \in \mathbb{Z}$$



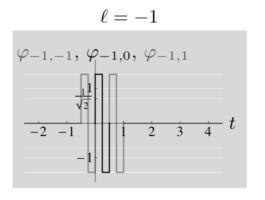


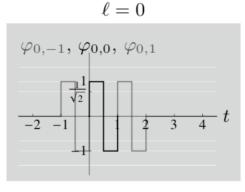
Piet Mondrian (1921)

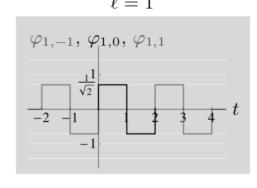


- example: consider a square wave function and $t_0=1$, a=2

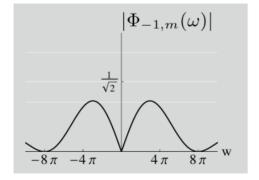
$$\varphi_{l,m}(t) = \varphi(\frac{t - 2^l m}{2^l}), \quad \varphi(t) = \begin{cases} 1, & \text{for } 0 \le t < \frac{1}{2}; \\ -1, & \text{for } \frac{1}{2} \le t < 1; \\ 0, & \text{otherwise.} \end{cases}$$

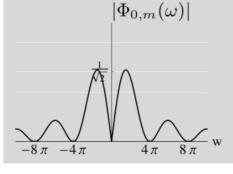


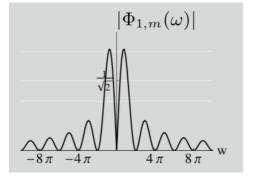




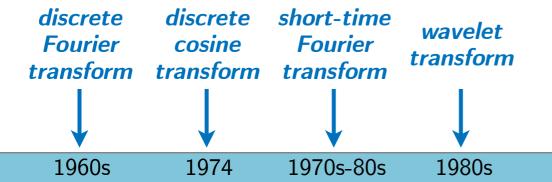
Basis functions.





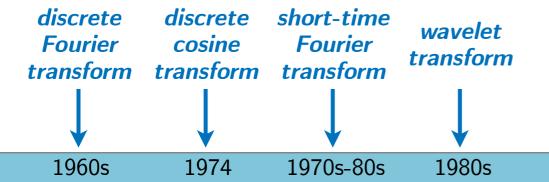


Magnitudes of the Fourier transform.



- consider a more general function and define the following transform

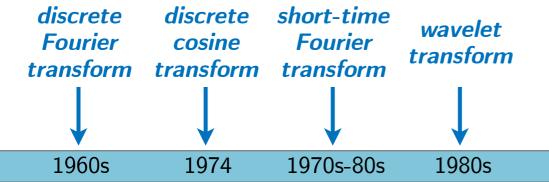
$$\psi_{s,\tau}(t) = \frac{1}{\sqrt{s}}\psi(\frac{t-\tau}{s}) \quad \Longrightarrow \quad X(s,\tau) = \int_{-\infty}^{\infty} x(t)\psi_{s,\tau}^*(t)dt = \int_{-\infty}^{\infty} x(t)\frac{1}{\sqrt{s}}\psi^*(\frac{t-\tau}{s})dt$$



- consider a more general function and define the following transform

$$\psi_{s,\tau}(t) = \frac{1}{\sqrt{s}} \psi(\frac{t-\tau}{s}) \longrightarrow X(s,\tau) = \int_{-\infty}^{\infty} x(t) \psi_{s,\tau}^*(t) dt = \int_{-\infty}^{\infty} x(t) \frac{1}{\sqrt{s}} \psi^*(\frac{t-\tau}{s}) dt$$

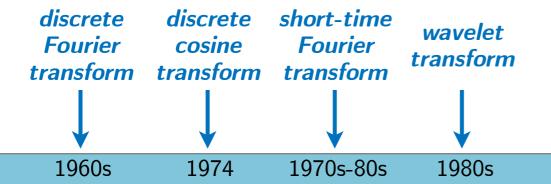
- the prototype function $\psi(t)$
 - has a compact support (small or "-let")
 - is band-pass with zero mean ("wave"): $\int_{-\infty}^{\infty} \psi(t) dt = 0$



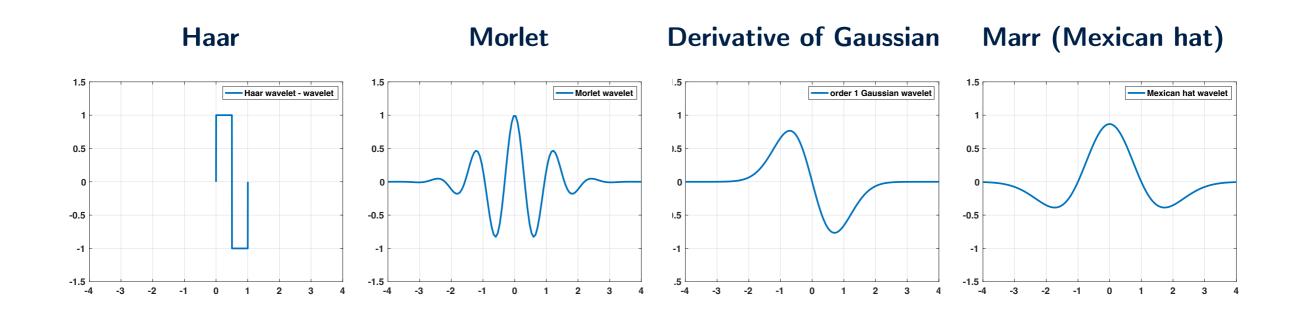
- consider a more general function and define the following transform

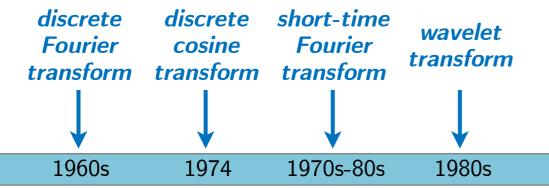
$$\psi_{s,\tau}(t) = \frac{1}{\sqrt{s}} \psi(\frac{t-\tau}{s}) \longrightarrow X(s,\tau) = \int_{-\infty}^{\infty} x(t) \psi_{s,\tau}^*(t) dt = \int_{-\infty}^{\infty} x(t) \frac{1}{\sqrt{s}} \psi^*(\frac{t-\tau}{s}) dt$$

- the prototype function $\psi(t)$
 - has a compact support (small or "-let")
 - is band-pass with zero mean ("wave"): $\int_{-\infty}^{\infty} \psi(t) dt = 0$
- this is called the continuous wavelet transform (CWT)
- CWT maps a 1-D function into a 2-D function (overcomplete)



examples of prototype function (mother wavelet)



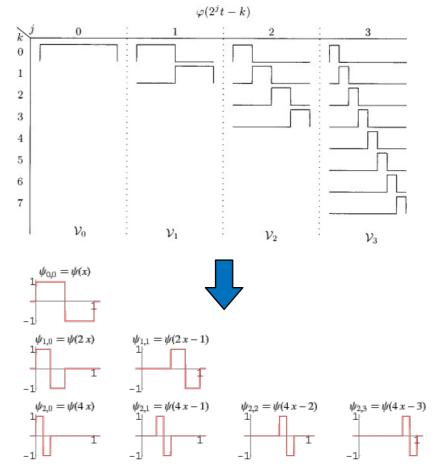


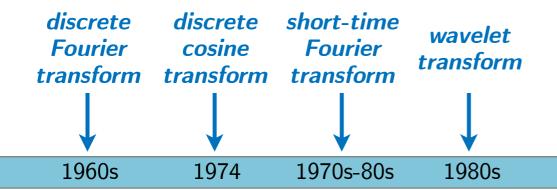
- CWT is an overcomplete transform; however, we can design an orthogonal wavelet transform through a multiresolution analysis

$$\psi_{l,m}(t) = \frac{1}{\sqrt{2l}}\psi(\frac{t-2^l m}{2^l})$$

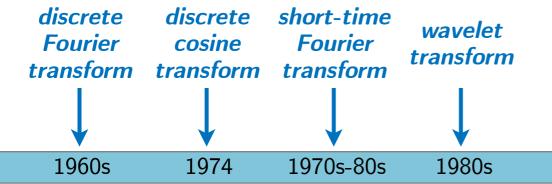
design principle

- form nested multiresolution spaces using scaling function φ
- obtain wavelets ψ by difference between nested subspaces

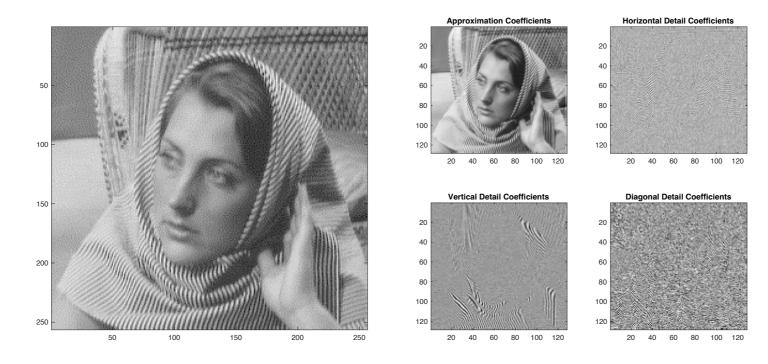


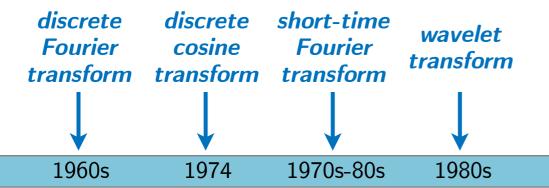


this leads to the discrete wavelet transform (DWT) which provides an orthogonal dictionary

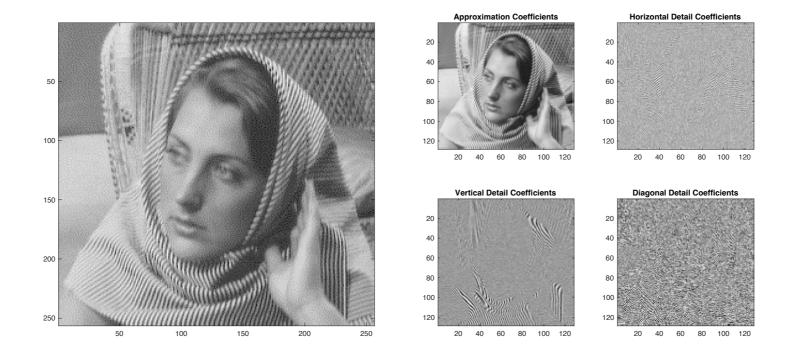


this leads to the discrete wavelet transform (DWT) which provides an orthogonal dictionary



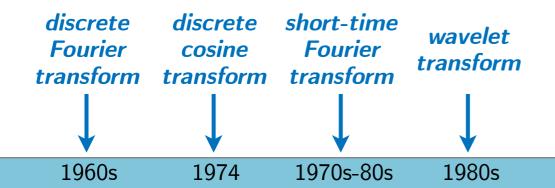


this leads to the discrete wavelet transform (DWT) which provides an orthogonal dictionary

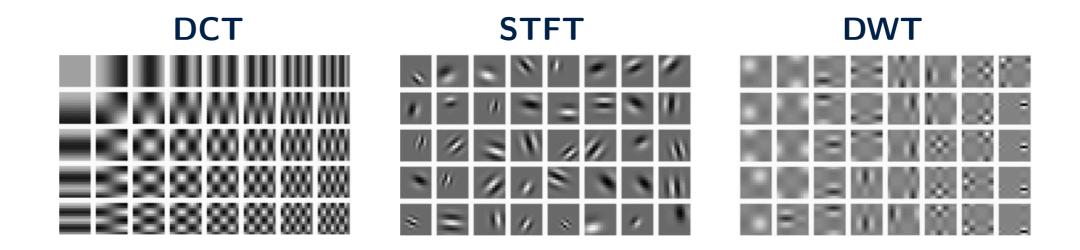


- DWT is behind the JEPG 2000 image compression standard

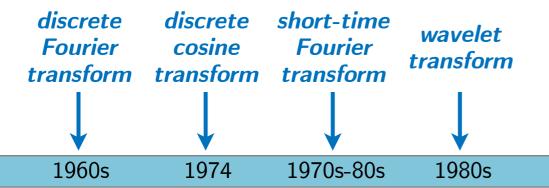
DCT vs STFT vs DWT



- comparison of the dictionaries we looked at so far

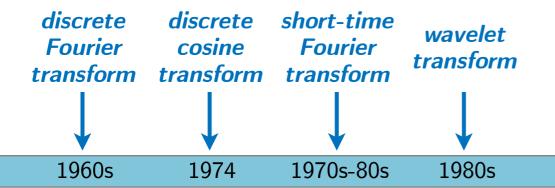


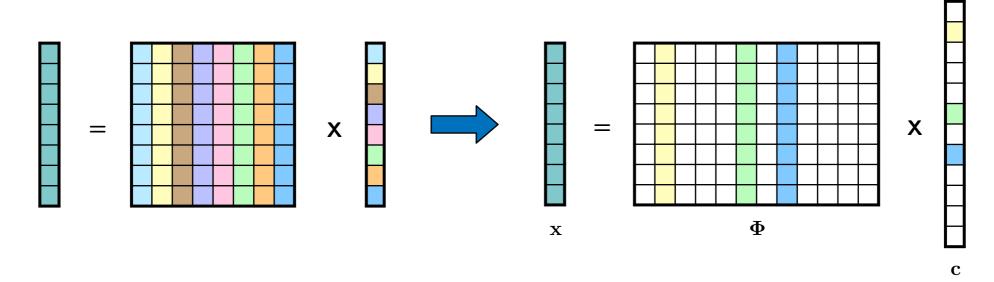
Transform/analytic dictionary design



- summary
 - modelling data by a simpler class of mathematical functions
 - smooth functions (DFT, DCT)
 - piecewise-smooth functions (wavelets)
 - desired properties
 - localisation (STFT, wavelets)
 - multiresolution (wavelets)
 - adaptivity (wavelet packets)
 - fast implementation is usually available
 - limited expressiveness

A paradigm shift in dictionary design





orthogonal atoms

complete dictionary

all signals use all atoms

dense coefficients

mathematical modelling

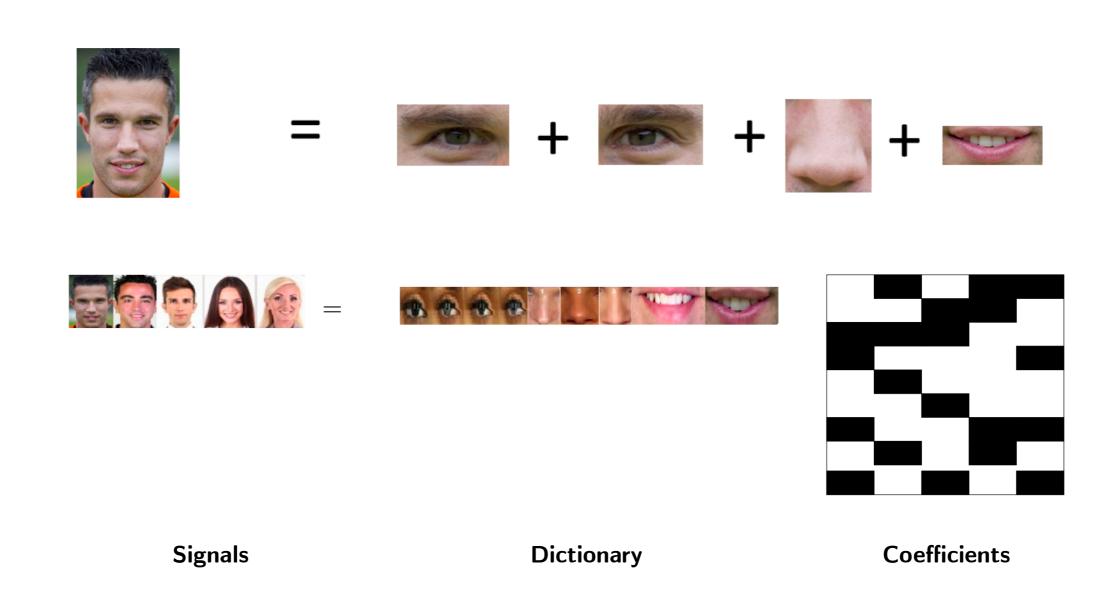
non-orthogonal atoms
overcomplete dictionary
different signals use different atoms
sparse coefficients
adaptation to data realisations

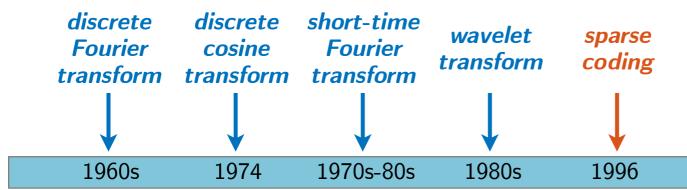
An illustrative example

• Modelling assumption: Each data point is a combination of only a few (sparse) fundamental elements, i.e., dictionary atoms

An illustrative example

• Modelling assumption: Each data point is a combination of only a few (sparse) fundamental elements, i.e., dictionary atoms

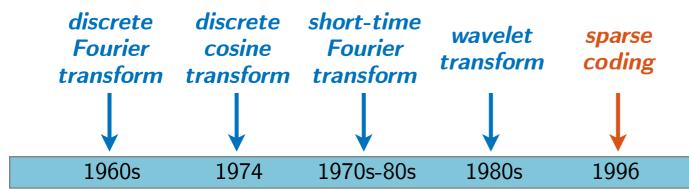




- the dictionary learning problem can be formulated as

$$\min_{\mathbf{\Phi}, \mathbf{c}} ||\mathbf{x} - \mathbf{\Phi} \mathbf{c}||_2^2 + \lambda ||\mathbf{c}||_1$$

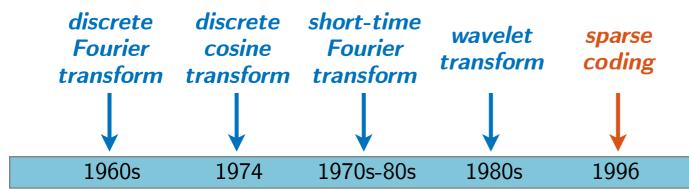
- **sparse approximation:** given Φ , solve for ${f c}$ via Lasso
- dictionary update: given c, update Φ via gradient descent



- the dictionary learning problem can be formulated as

$$\min_{\mathbf{\Phi}, \mathbf{c}} ||\mathbf{x} - \mathbf{\Phi} \mathbf{c}||_2^2 + \lambda ||\mathbf{c}||_1$$

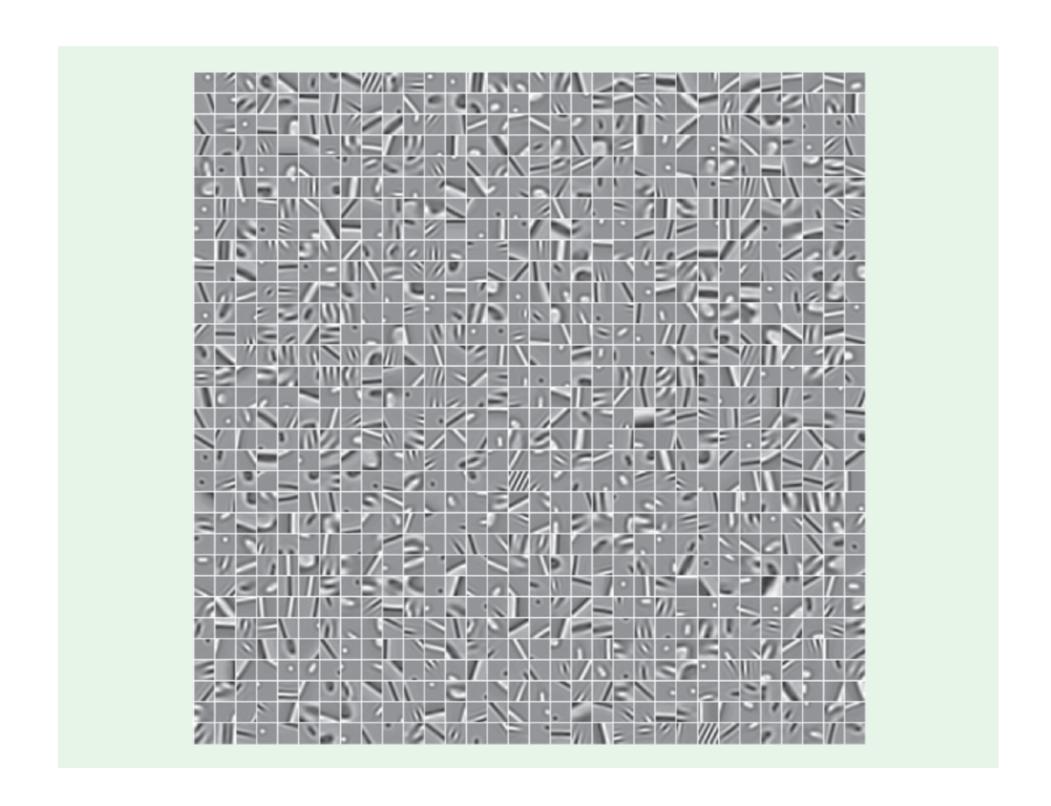
- **sparse approximation:** given Φ , solve for ${f c}$ via Lasso
- ullet dictionary update: given $oldsymbol{c}$, update $oldsymbol{\Phi}$ via gradient descent
- works at patch level for efficiency
- does not necessarily find global optimum

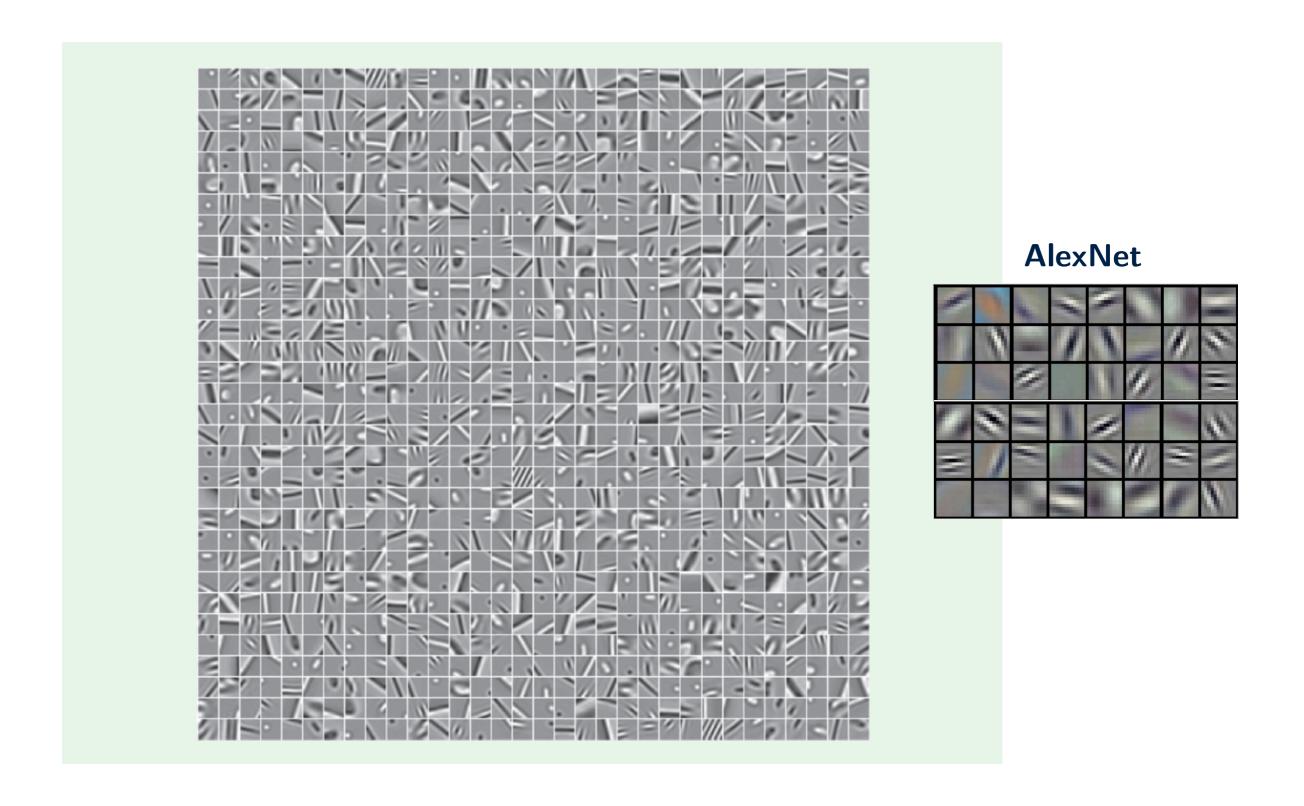


- the dictionary learning problem can be formulated as

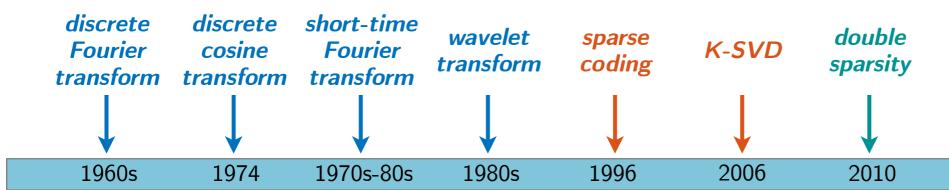
$$\min_{\mathbf{\Phi},\mathbf{c}} ||\mathbf{x} - \mathbf{\Phi}\mathbf{c}||_2^2 + \lambda ||\mathbf{c}||_1$$

- sparse approximation: given Φ , solve for ${f c}$ via Lasso
- dictionary update: given ${f c}$, update ${f \Phi}$ via gradient descent
- works at patch level for efficiency
- does not necessarily find global optimum
- trained atoms are remarkably similar to mammalian simple-cell receptive fields





Dictionary learning

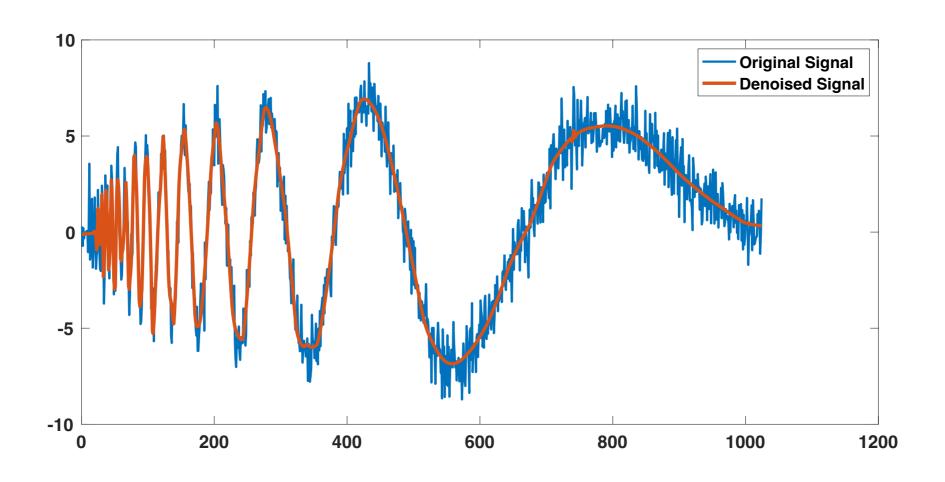


- summary
 - learning representations directly from data realisations
 - desired properties
 - overcompleteness
 - sparse representations
 - efficiency in training
 - may be combined with analytical dictionary design
 - trained dictionary with structures (e.g., parametric dictionary learning via double sparsity)

Lecture 1

- Introduction & Basic concepts and tools
- A historical overview of signal representation techniques
- Applications & Discussion

Application I: Signal denoising



denoising using the order 4 symlet wavelets

Application II: Image compression

original



JPEG 2000 (10% in size)



JPEG 2000 (1% in size)



compression using the Cohen-Daubechies-Feauveau wavelets

Application III: Image reconstruction

50 % missing pixels



Average # coeffs: 4.0202 Average # coeffs: 4.7677 Average # coeffs: 4.7694 MAE: 0.012977 RMSE: 0.029204

Learned reconstruction



Learned reconstruction MAE: 0.020035 RMSE: 0.055643





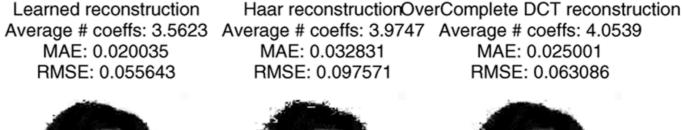
MAE: 0.025001

MAE: 0.015719

RMSE: 0.037745

Haar reconstructionOverComplete DCT reconstruction





MAE: 0.022833

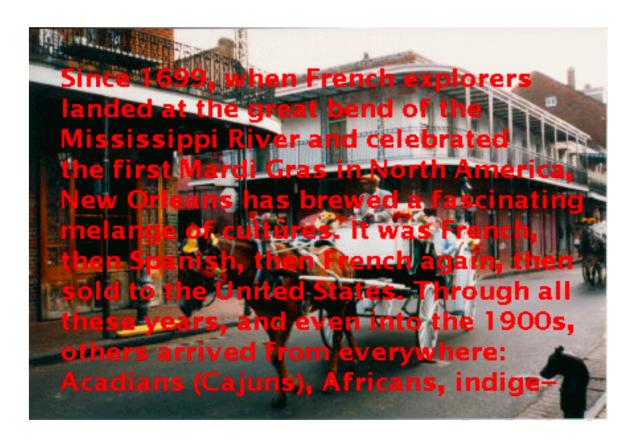
RMSE: 0.071107



70 % missing pixels

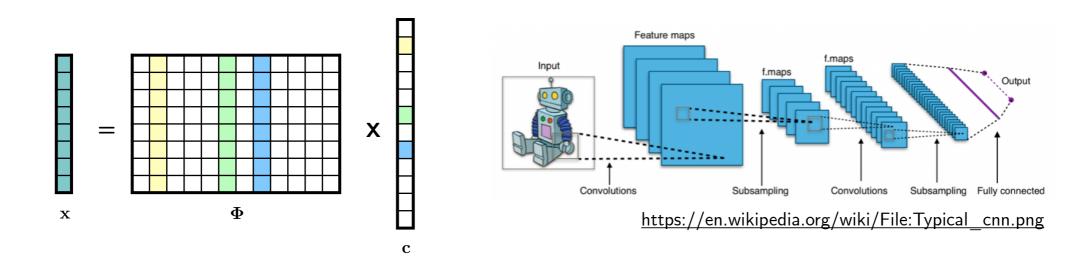


Application IV: Image restoration

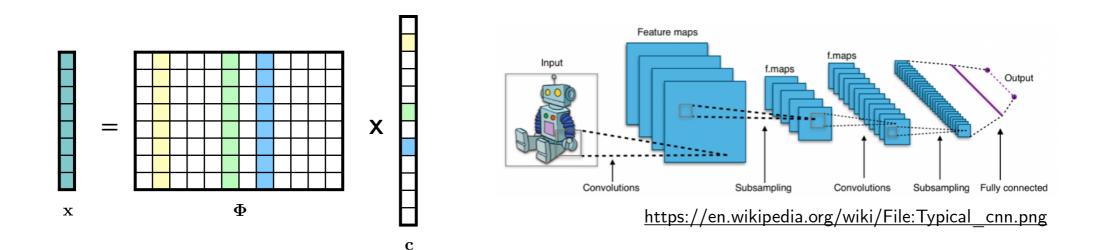




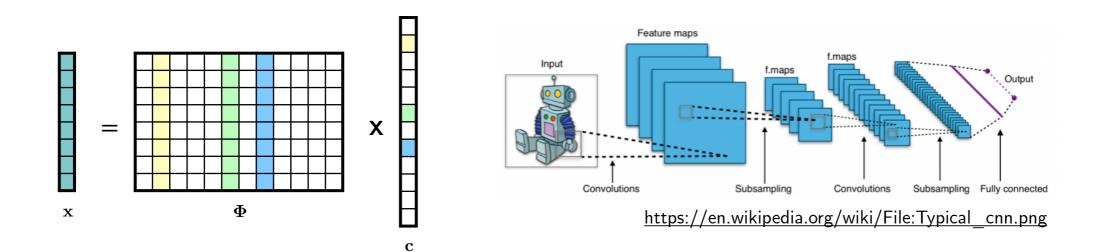
Dictionary learning vs Deep learning



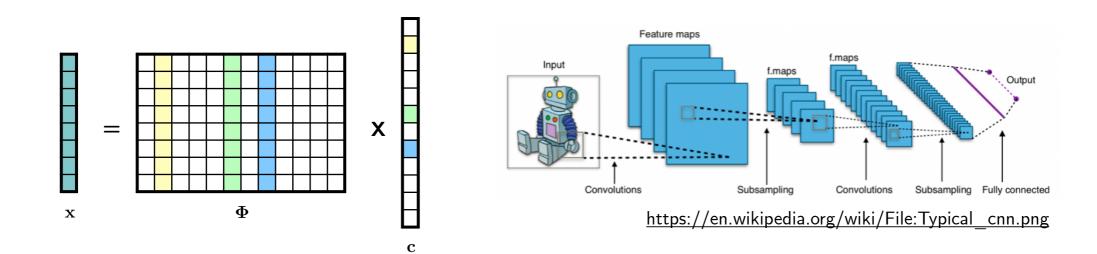
- Dictionary learning vs Deep learning
 - both extract feature representations from data realisations



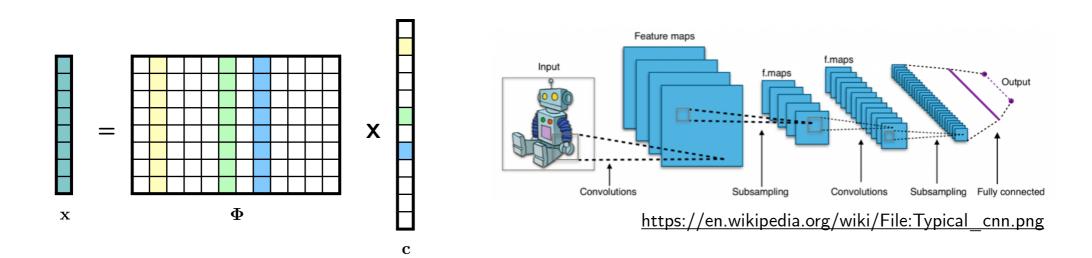
- Dictionary learning vs Deep learning
 - both extract **feature representations** from data realisations
 - both apply sparsifying operations such as shrinkage or rectified linear units



- Dictionary learning vs Deep learning
 - both extract feature representations from data realisations
 - both apply sparsifying operations such as shrinkage or rectified linear units
 - the former often leads to **shallow** representations while the latter to **hierarchical** representations (via multi-layered convolution and pooling)

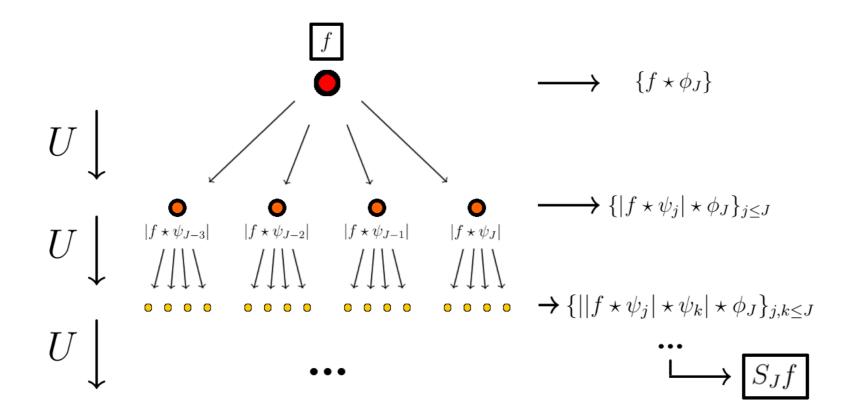


- Dictionary learning vs Deep learning
 - both extract feature representations from data realisations
 - both apply sparsifying operations such as shrinkage or rectified linear units
 - the former often leads to **shallow** representations while the latter to **hierarchical** representations (via multi-layered convolution and pooling)
 - the former is mainly for **reconstruction/approximation** (similar to autoencoders) while the latter is widely used for **classification**



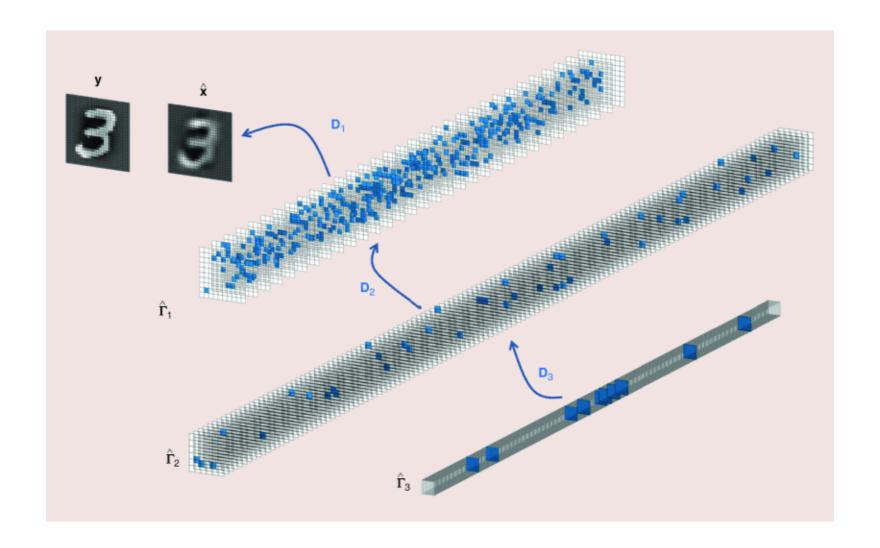
Dictionary-inspired deep architectures

Scattering transform



Dictionary-inspired deep architectures

Multi-layer convolutional sparse coding



References



Dictionaries for Sparse Representation Modeling

Digital sampling can display signals, and it should be possible to expose a large part of the desired signal information with only a limited signal sample.

By Ron Rubinstein, Student Member IEEE, Alfred M. Bruckstein, Member IEEE, and MICHAEL ELAD, Senior Member IEEE

ABSTRACT | Sparse and redundant representation modeling of data assumes an ability to describe signals as linear combinations of a few atoms from a pre-specified dictionary. As such, signal processing techniques commonly require the choice of the dictionary that sparsifies the signals is crucial more meaningful representations which capture the useful for the success of this model. In general, the choice of a proper characteristics of the signal—for recognition, the repredictionary can be done using one of two ways: i) building a sparsifying dictionary based on a mathematical model of the data, or ii) learning a dictionary to perform best on a training noise; and for compression, the representation should set. In this paper we describe the evolution of these two capture a large part of the signal with only a few coefficients. Interestingly, in many cases these seemingly differences such as wavelets, wavelet packets, contourlets, and curvelets, all aiming to exploit 1-0 and 2-D mathematical modes for constructing effective dictionaries for signals and images.

Dictionary learning takes a different route, attaching the decompose the signal. When the dictionary forms a basis, dictionary to a set of examples it is supposed to serve. From the seminal work of Field and Olshausen, through the MOD, the nation of the dictionary atoms. In the simplest case the K-SVD, the Generalized PCA and others, this paper surveys the dictionary is orthogonal, and the representation coeffi

KEYWORDS | Dictionary learning: harmonic analysis; signal also referred to as the bi-orthogonal dictionary approximation; signal representation; sparse coding; sparse

For years, orthogonal and bi-orthogonal

I. INTRODUCTION

The process of digitally sampling a natural signal leads to its

Manuscript received April 5, 2009; accepted November 21, 2009. Date of publication April 22, 2010; date of current version May 19, 2010. This remarch was pathly assigned by the European Community PFP-PET program, MALL project, under gont agreement 25913, and by the 15° grant 599/08. MLL project, under the subness are with the Department of Computer Science, The Technico-stude institute of Technology, Isa'd 33000, lorsal (e-mail: on multiplicated-inion-at-like throdity-cated-inion-at-like indicated-inion-at-like indicated-inion-a

0018-9219/\$26.00 @2010 IEEE

various options such training has to offer, up to the most recent cients can be computed as inner products of the signal and the atoms; in the non-orthogonal case, the cofficients are the inner products of the signal and the dictionary inverse,

For years, orthogonal and bi-orthogonal dictionaries were dominant due to their mathematical simplicity. How ever, the weakness of these dictionaries—namely their limited expressiveness—eventually outweighed their simplicity. This led to the development of newer overcomplete representation as the sum of Delta functions in space or signal, which promised to represent a wider range of signal

> The move to overcomplete dictionaries was done cau tiously, in an attempt to minimize the loss of favorable properties offered by orthogonal transforms. Many dictio-naries formed tight frames, which ensured that the repre-sentation of the signal as a linear combination of the atoms could still be identified with the inner products of the signal and the dictionary. Another approach, manifested by

Vol. 98, No. 6, June 2010 | PROCEEDINGS OF THE IEEE 1045

Wana Tošić and Pascal Frossard **Dictionary Learning** What is the right representation for my signal? uge amounts of high-dimensional information are captured every second by diverse natural sensors such as the eyes or ears, as well as artificial sensor like cameras or microphones. This information is largely redundant in two main aspects; it physical world and each version is usually densely sampled by generic sensors. The relevant information about the underlying processes that cause our observations is generally of much reduced dimensionality compared to such recorded data sets. The extraction of this relevant information by identifying the general ing causes within classes of signals is the central topic of this article. We present methods for determining the proper representation of data sets by means of tality subspaces, which are adaptive to both the characteristics of the signals and the processing task at hand. These representations are based on the principle that our observations can be described by a sparse subset of atoms taken from a redundant dictionary, which represents the causes of our observations of the world. We describe methods for learning dictionaries that are appropriate for the representation of given classes of signals and multisensor data. We further show that dim reduction based on dictionary representation can be extended to address specific tasks such as data analy sis or classification when the learning includes a class separability criteria in the objective function. The benefits of dictionary learning clearly show that a proper understanding of causes underlying the sensed world is key to task-specific representation of relevant information in high-dimensional data sets.

WHAT IS THE GOAL OF DIMENSIONALITY REDUCTION?

Natural and artificial sensors are the only tools we have for sensing the world and gathering information about physical processes and their causes. These sensors are usually not aware of the physical process underlying the phenomena they "see," hence they often sample the information with a higher rate than the effective dimension of the process. However, to store, transmit or analyze the processes we observe we do not need such abundant data: we only need the information that is relevant to understand the causes, to reproduce the physical processes, or to make decisions. In other words, we can reduce the

Digital Object Identifier 10.1109/MSP.2010.53953

IEEE SIGNAL PROCESSING MAGAZINE [27] MARCH 2011

1053-5888/11/626.006/2011/EEE

IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE. VOL. 35. NO. 8. AUGUST 2013

Representation Learning: A Review and New Perspectives

Yoshua Bengio, Aaron Courville, and Pascal Vincent

Abstract—The success of machine learning algorithms generally depends on data representation, and we hypothesize that this is because different representations can entangle and hide more or less the different explanatory factors of variation behind the data. Although specific domain knowledge can be used to help design prepresentations, learning with generic priors can also be used, and the quest for Al is motivating the design of more powerful representation-teaming algorithms implementing such priors. This paper reviews recent work in the area of unsupervised feature learning and deep learning, covering advances in probabilistic models, calencoders, manifold learning, and deep networks. This motivates longer term unanswered questions about the appropriate objectives for learning good representations, for computing representations (i.e., inference), and the geometrical connections between representation learning, density estimation, and manifold learning.

Index Terms—Deep learning, representation learning, feature learning, unsupervised learning, Boltzmann machine, autoencoder neural nets

1 INTRODUCTION

data transformations that result in a representation of the data that can support effective machine learning. Such feature engineering is important but labor intensive and highlights the weakness of current learning algorithms. Their inability to extract and organize the discriminative information from the data. Feature engineering is a way to take advantage of human ingenuity and prior knowledge to compensate for that weakness. To expand the scope and ease of applicability of machine learning, it would be highly desirable to make learning algorithms less dependent on feature engineering so that novel applications could be constructed faster, and more importantly, to make progress toward artificial intelligence (AI). An Al must fundamental transformations with the goal of yielding more abstract—and ultimately more useful—representations. Here, we survey this rapidly developing area with special emphasis on recent progress. We even driving research in this area. Specifically, what makes one representation better than another? Given an example how should we compute its representation, i.e., perform feature engineering so that novel applications could be constructed faster, and more importantly, to make progress to a survey of the progression of multiple nonlinear transformations with the goal of yielding more abstract—and ultimately more useful—representations. Here, we survey this rapidly developing area with special emphasis on recent progress. Weet of the fundamental questions that have been driving research in this area. Specifically, what makes or representation better than another? Given an example how should we compute its representation. Here, we survey this rapidly developing area with special emphasis on recent progress. Here, we survey this rapidly developing area with special emphasis on recent progress. Here, we survey this rapidly developing area with special emphasis on recent progress. Here, we survey this rapidly developing area with special emphasis on recent progress. Here, we survey this rapid toward artificial intelligence (AI). An AI must fundamentally understand the world around us, and we argue that this can only be achieved if it can learn to identify and

the observed milieu of low-level sensory data.

This paper is about representation learning, i.e., learning representations of the data that make it easier to extract In spaper is about representation learning, i.e., learning representations of the data that make it easier to extract header of Deep Learning or Feature Learning, Although depth useful information when building classifiers or other

numscript received 9 Apr. 2012; received 17 Oct. 2012; accepted 24 Feb. 2013; highlight some of these high points. Whished enline 28 Feb. 2013. Post of the points of the points of the points of the points. Speech Recognition and Signal Processing

Recommended for acceptance by S. Bengio, L. Deng, H. Larochelle, H. Lee, and
S. Salakhutdinov.
S. Salakhutdinov.
See Salakhutdinov.
Speech mecognition and Signal Processing
Speech was one of the early applications of neural networks, in particular convolutional (or time-delay) neural TPAMIS-201-204-260.

Digital Object Identifier no. 10.1109/TPAMI.2013.50.

0162-8828/13/\$31.00 © 2013 IEEE Published by the IEEE Computer Society

The performance of machine learning methods is heavily dependent on the choice of data representation (or features) on which they are applied. For that reason, much of the actual effort in deploying machine learning algorithms goes into the design of preprocessing pipelines and data transformations that result in a representation of the the data transformations that result in a representation of the theorem of the underlying explanatory factors for the observed input. A good representation is also one that is useful in a proper in the composition of the underlying explanatory factors for the observed input. A good representation is also one that is useful in a input to a supervised prefer in the properties of the presentation of the underlying explanatory factors for the observed input. A good representation is also one that is useful in the properties of the prop

REPRESENTATIONS?

Representation learning has become a field in itself in the can only be achieved if it can learn to identify and disentangle the underlying explanatory factors hidden in the observed milieu of low-level sensory data.

**Representation learning has become a field in the learning community, with regular workshops at the leading conferences such as NIPS and ICML, and a new conference dedicated to it. ICLR.1 sometimes under the predictors. In the case of probabilistic models, a good interesting and can be conveniently captured when the representation is often one that captures the posterior problem is cast as one of learning a representation, as discussed in the next section. The rapid increase in scientific activity on representation learning has been accompanied The authors are with the Department of Computer Science and Operations Research, Université de Montréal, Do Box 6128, Succ. Centre-Ville, Montreal, Quebe 1H2 37, Cranda.
 Des constitution of the State of Computer Science and Operations and nourished by a remarkable string of empirical successes both in academia and in industry. Below, we briefly

1. International Conference on Learning Representations

66/66