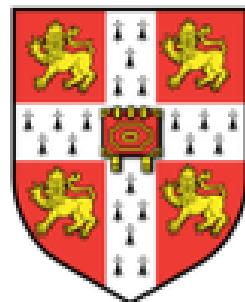


Gaussian Processes for Audio Feature Extraction

Dr. Richard E. Turner (ret26@cam.ac.uk)

Computational and Biological Learning Lab
Department of Engineering
University of Cambridge



UNIVERSITY OF
CAMBRIDGE

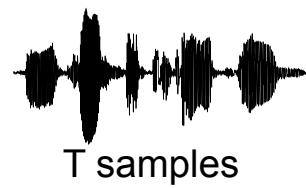


Computational and
Biological Learning

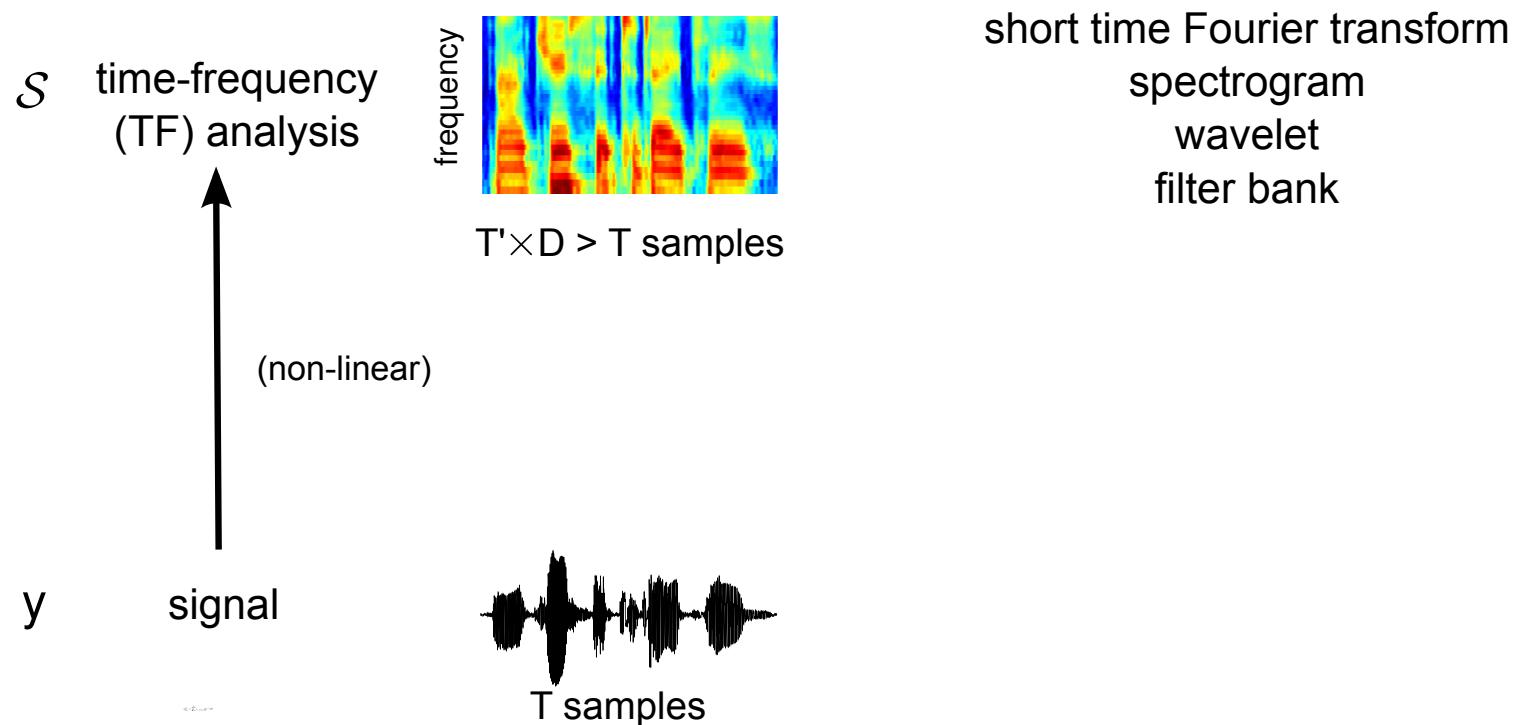
Machine hearing pipeline

y

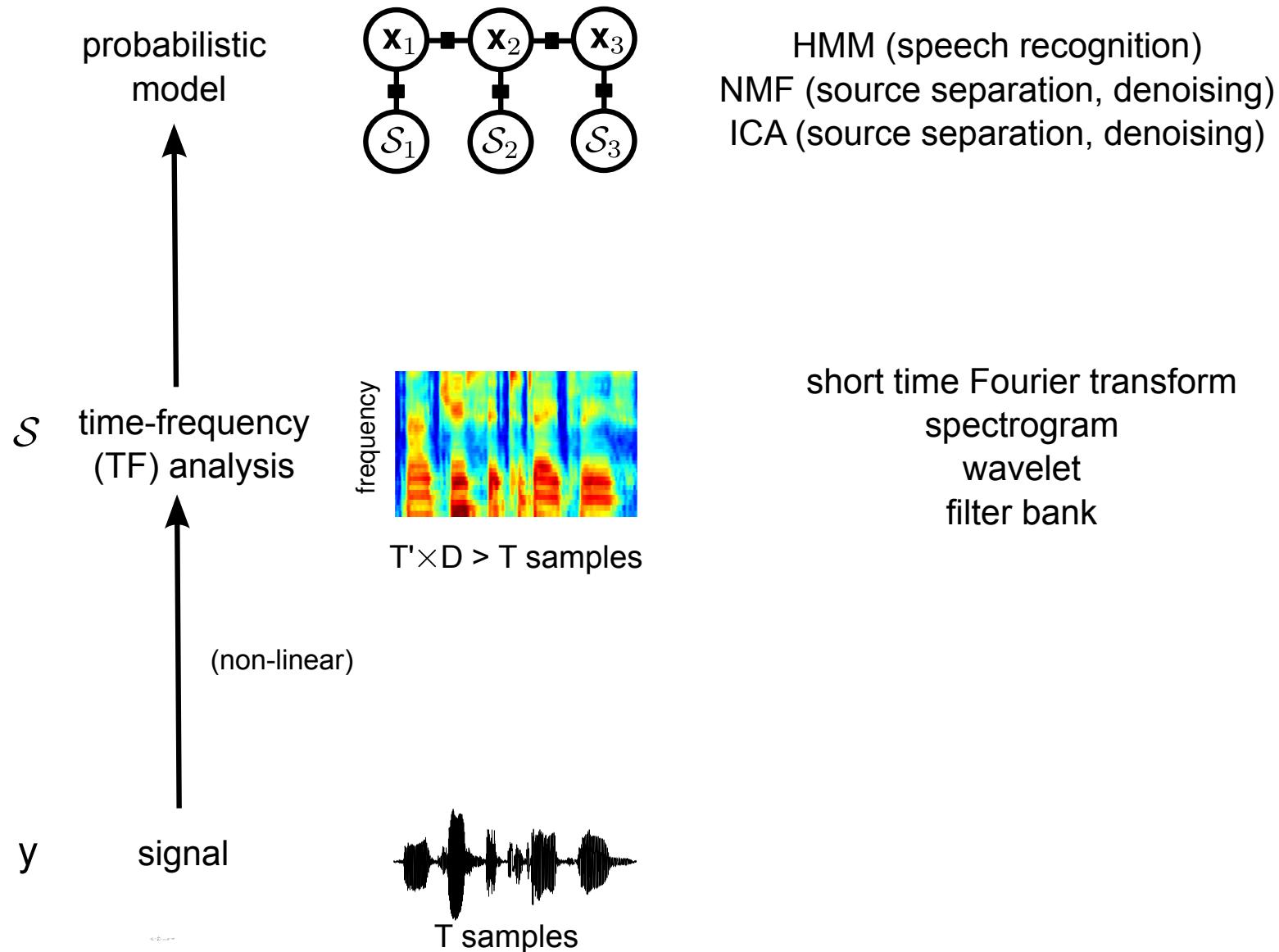
signal



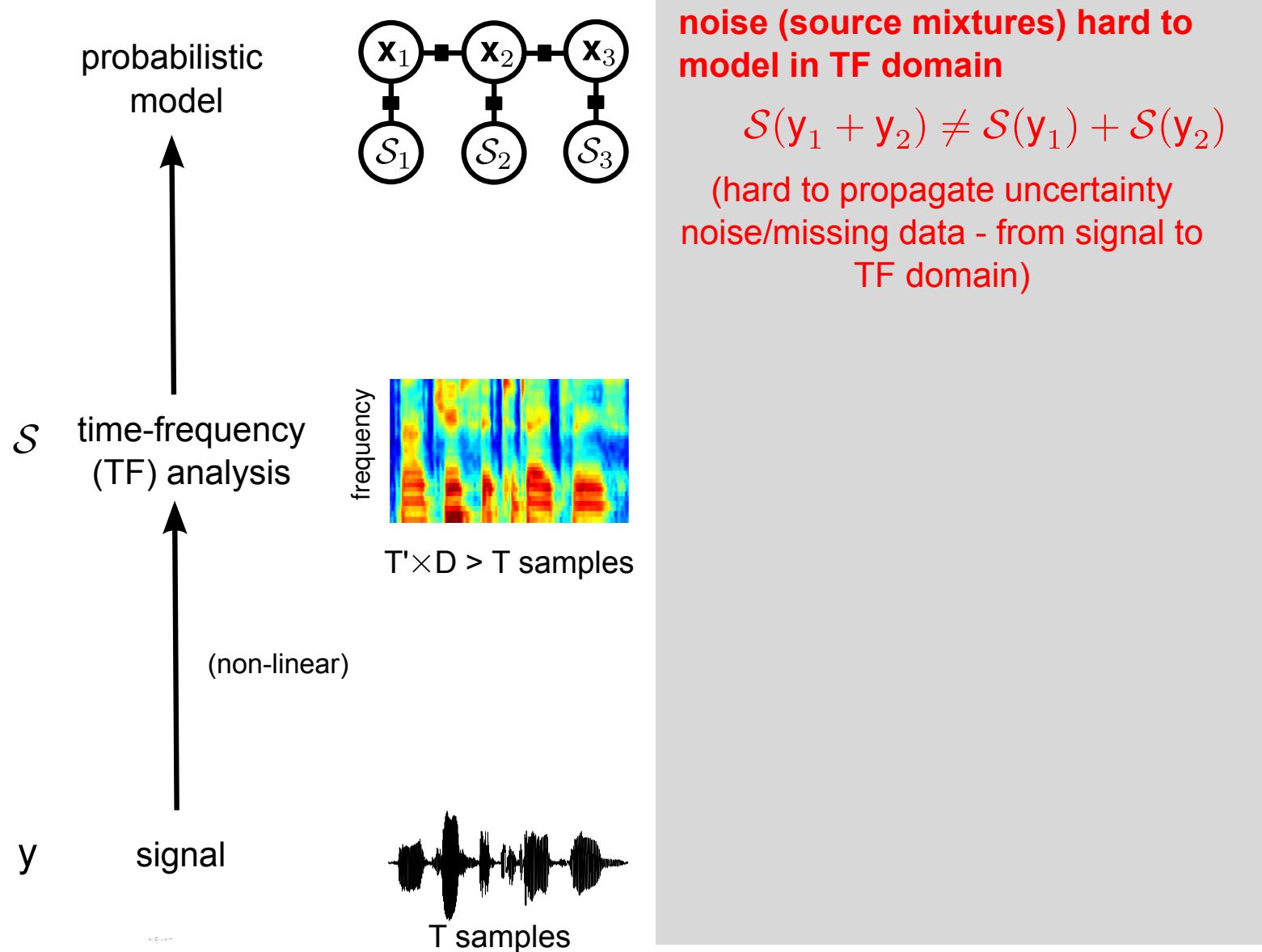
Machine hearing pipeline



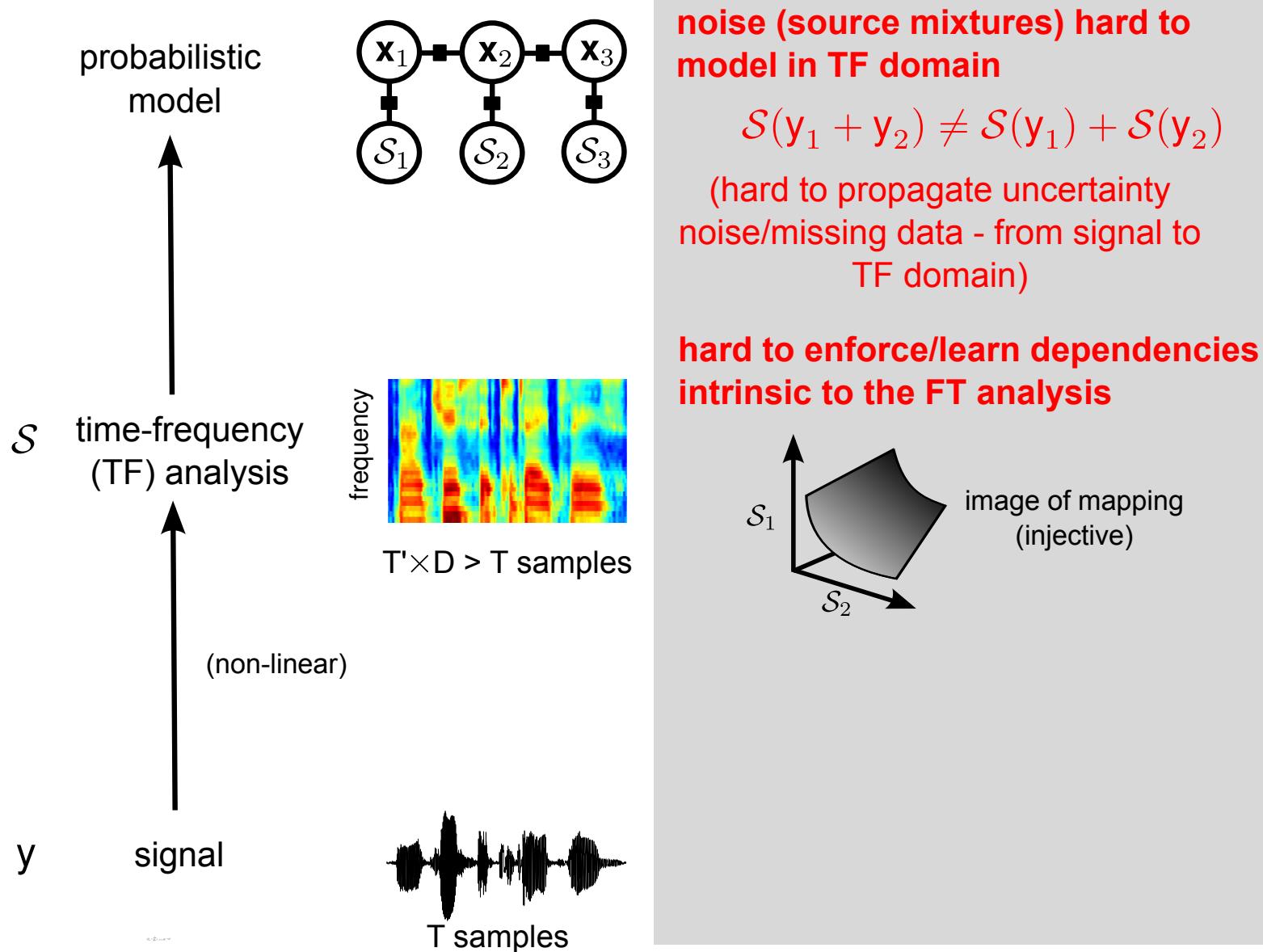
Machine hearing pipeline



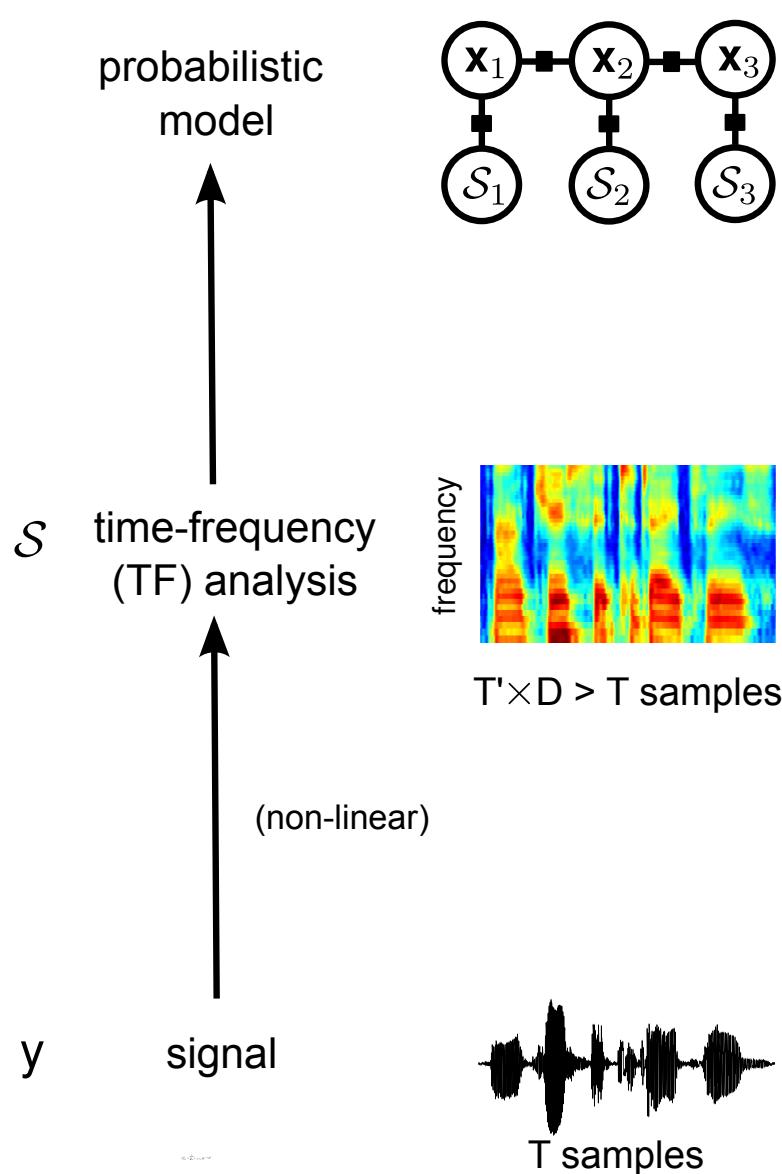
Problems with conventional pipeline



Problems with conventional pipeline



Problems with conventional pipeline

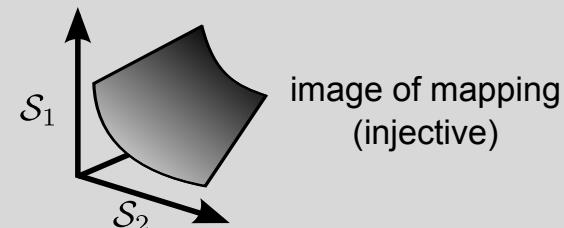


noise (source mixtures) hard to model in TF domain

$$\mathcal{S}(y_1 + y_2) \neq \mathcal{S}(y_1) + \mathcal{S}(y_2)$$

(hard to propagate uncertainty noise/missing data - from signal to TF domain)

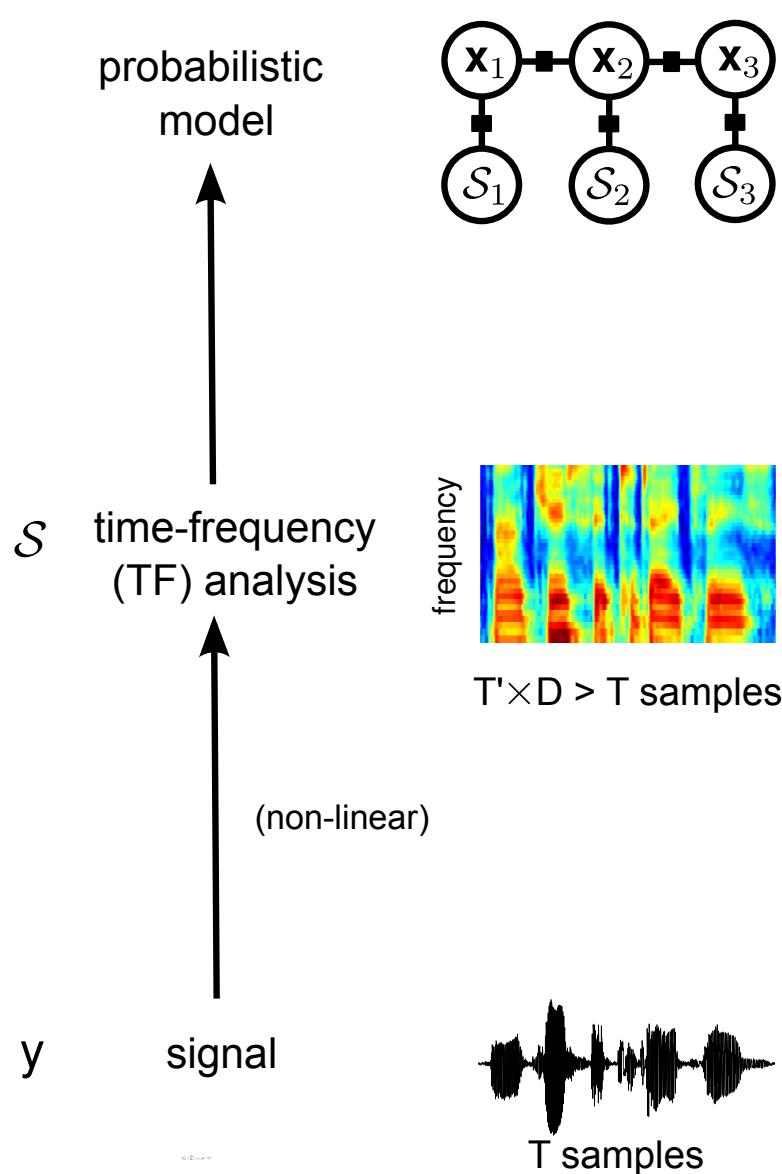
hard to enforce/learn dependencies intrinsic to the FT analysis



learning based on time-frequency representation ignores Jacobian

$$p(\mathcal{S}|\theta)$$

Problems with conventional pipeline



noise (source mixtures) hard to model in TF domain

$$\mathcal{S}(y_1 + y_2) \neq \mathcal{S}(y_1) + \mathcal{S}(y_2)$$

(hard to propagate uncertainty noise/missing data - from signal to TF domain)

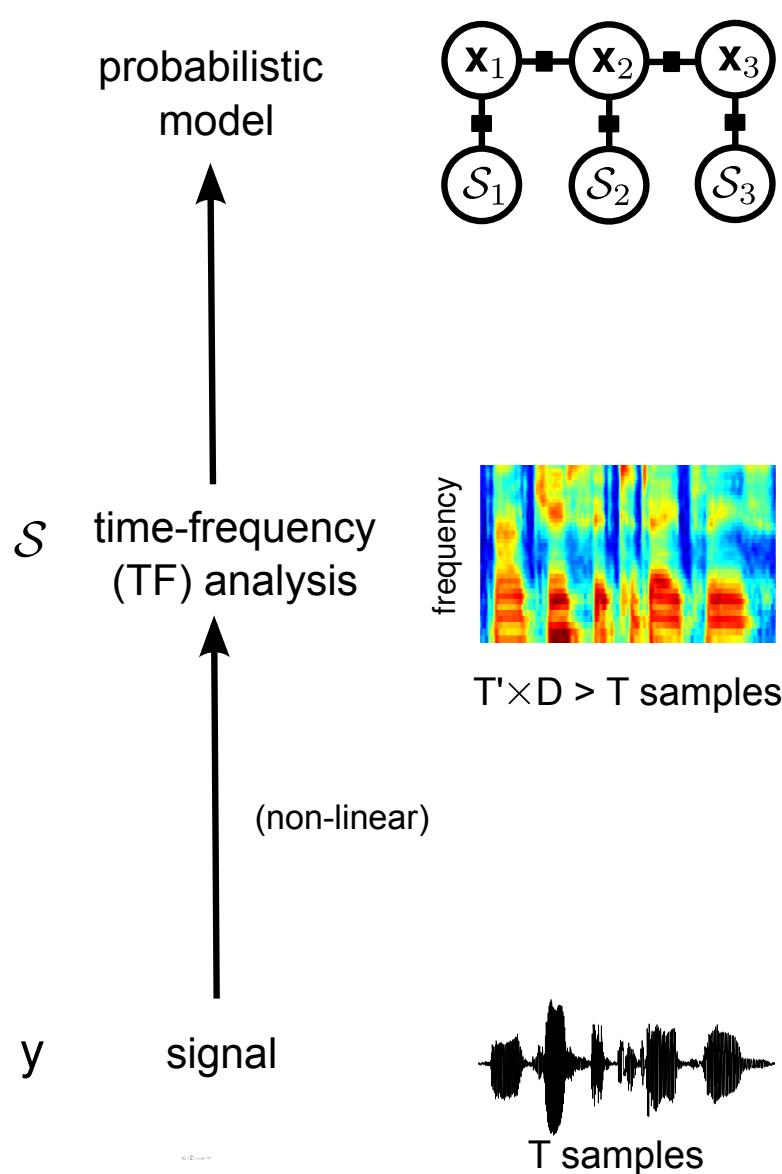
hard to enforce/learn dependencies intrinsic to the FT analysis

A diagram showing a shaded, curved surface representing the image of a mapping from \mathcal{S}_1 to \mathcal{S}_2 . The axes are labeled \mathcal{S}_1 and \mathcal{S}_2 . The text "image of mapping (injective)" is written next to the surface.

learning based on time-frequency representation ignores Jacobian

$$p(y|\theta) = p(\mathcal{S}|\theta) \left| \frac{d\mathcal{S}}{dy} \right|$$

Problems with conventional pipeline



noise (source mixtures) hard to model in TF domain

$$\mathcal{S}(y_1 + y_2) \neq \mathcal{S}(y_1) + \mathcal{S}(y_2)$$

(hard to propagate uncertainty noise/missing data - from signal to TF domain)

hard to enforce/learn dependencies intrinsic to the FT analysis

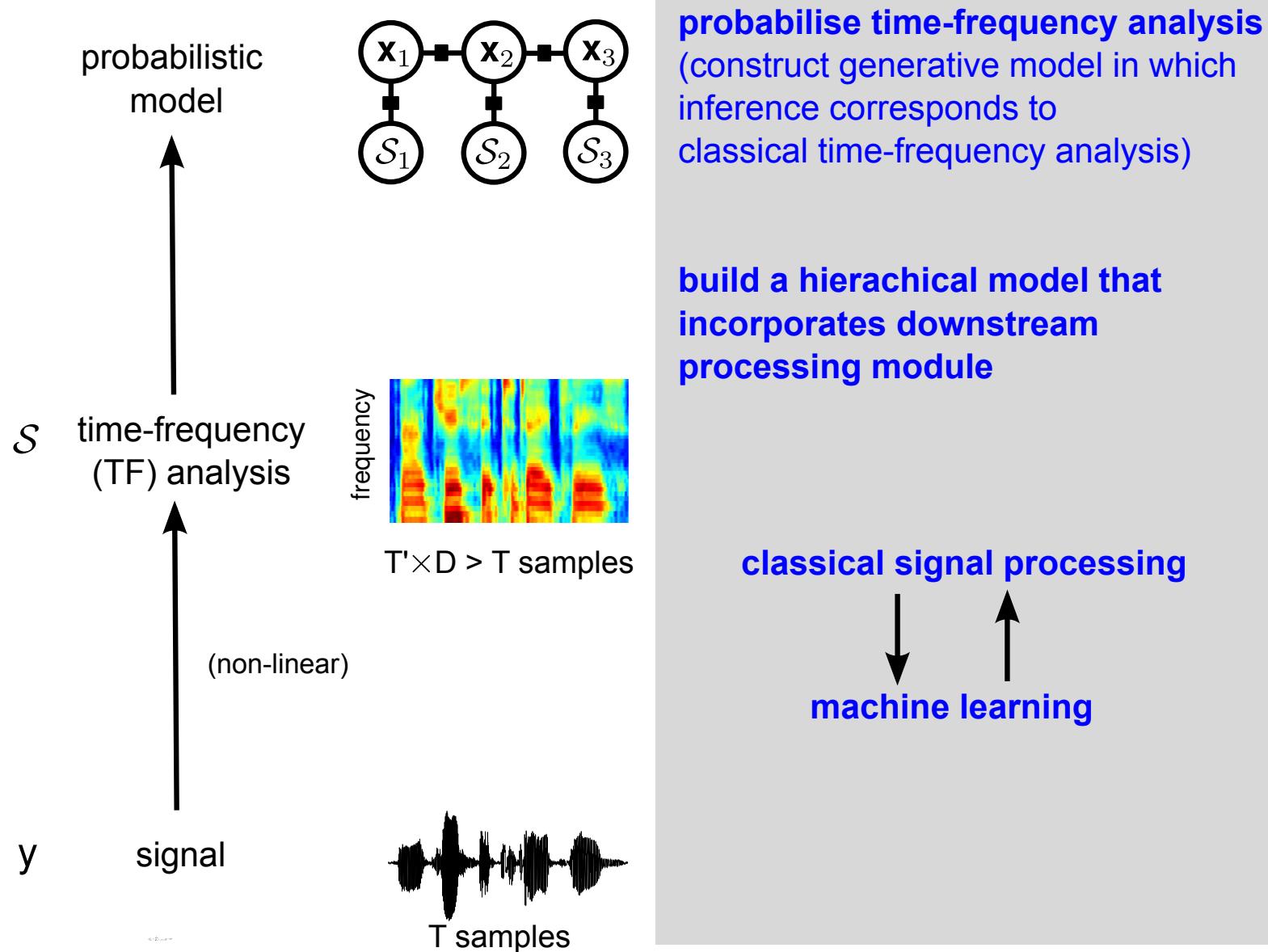
A diagram showing a shaded, curved surface representing the "image of mapping (injective)" from a space S_1 to a space S_2 . The mapping is injective, meaning different points in S_1 map to different points in S_2 .

learning based on time-frequency representation ignores Jacobian

$$p(y|\theta) = p(\mathcal{S}|\theta) \left| \frac{d\mathcal{S}}{dy} \right|$$

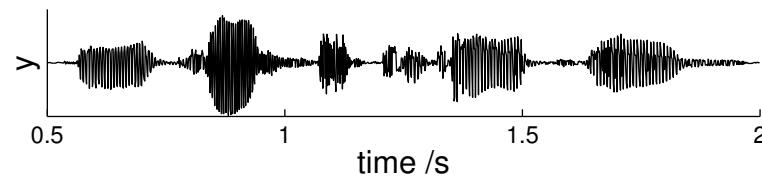
hard to adapt both top and bottom layers

Goal of this talk



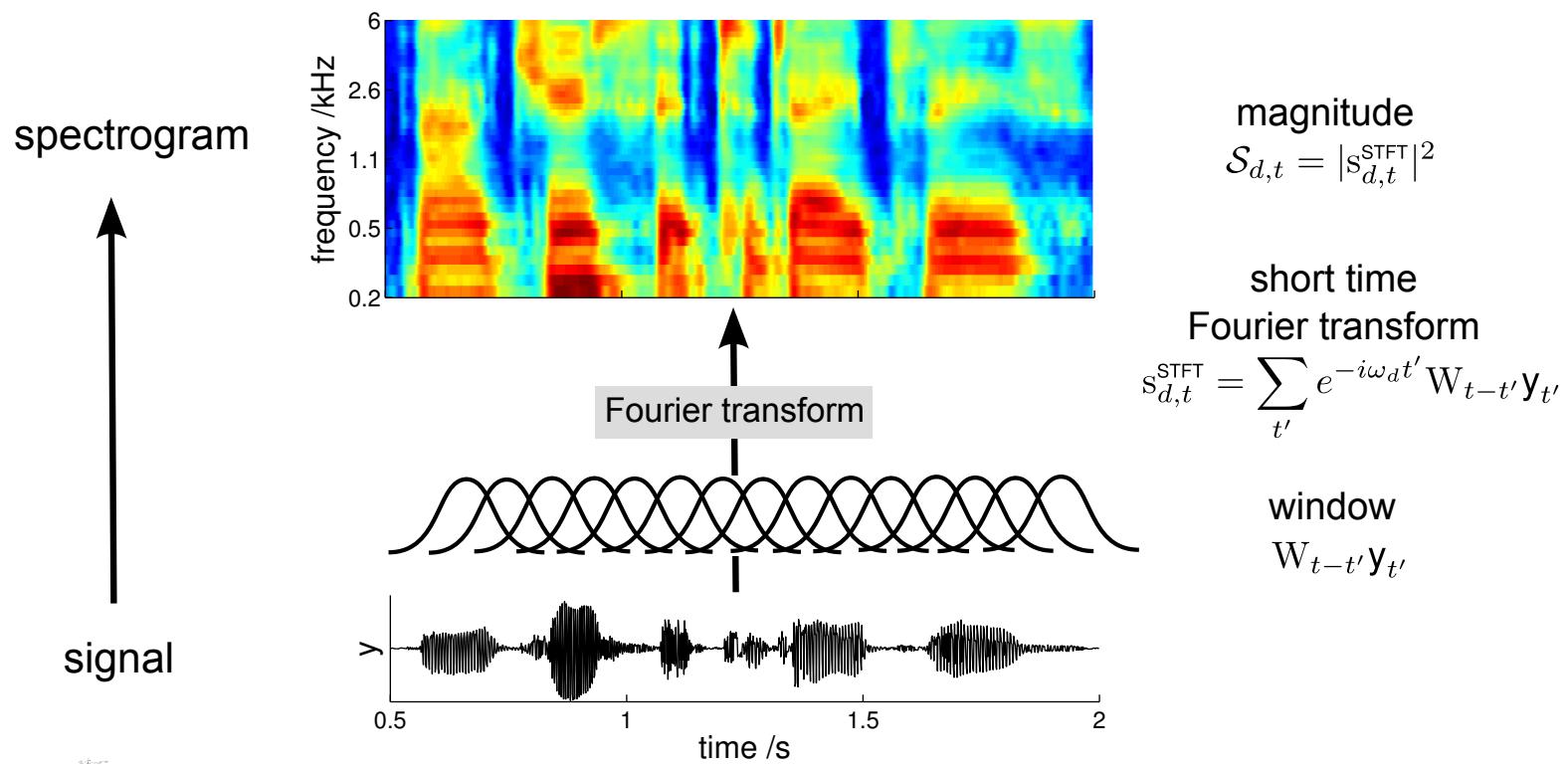
A typical audio pipeline

signal

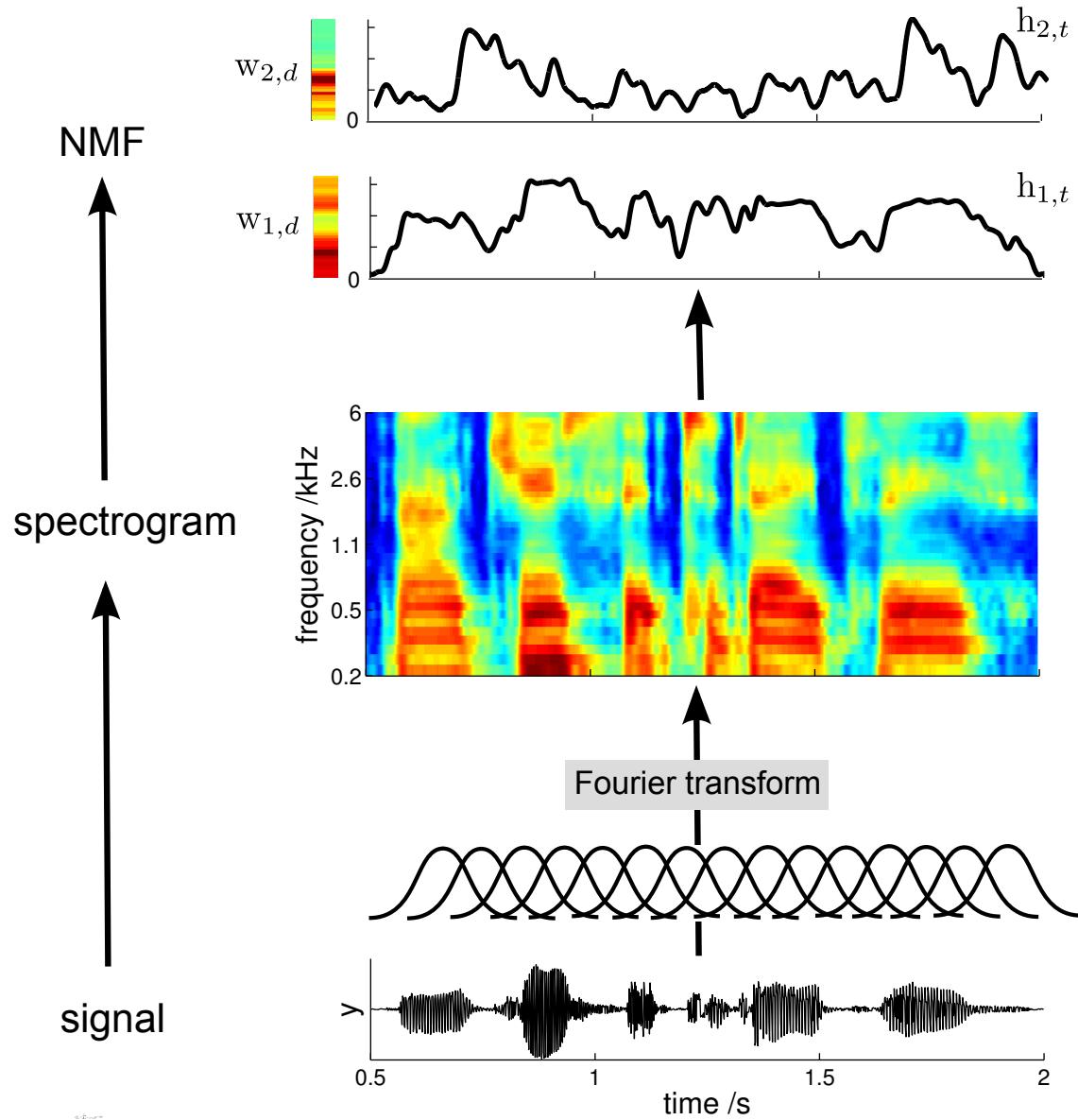


© R. C. Miller

A typical audio pipeline



A typical audio pipeline



$$s_{d,t} = \sum_{l=1}^L h_{l,t} w_{l,d}$$

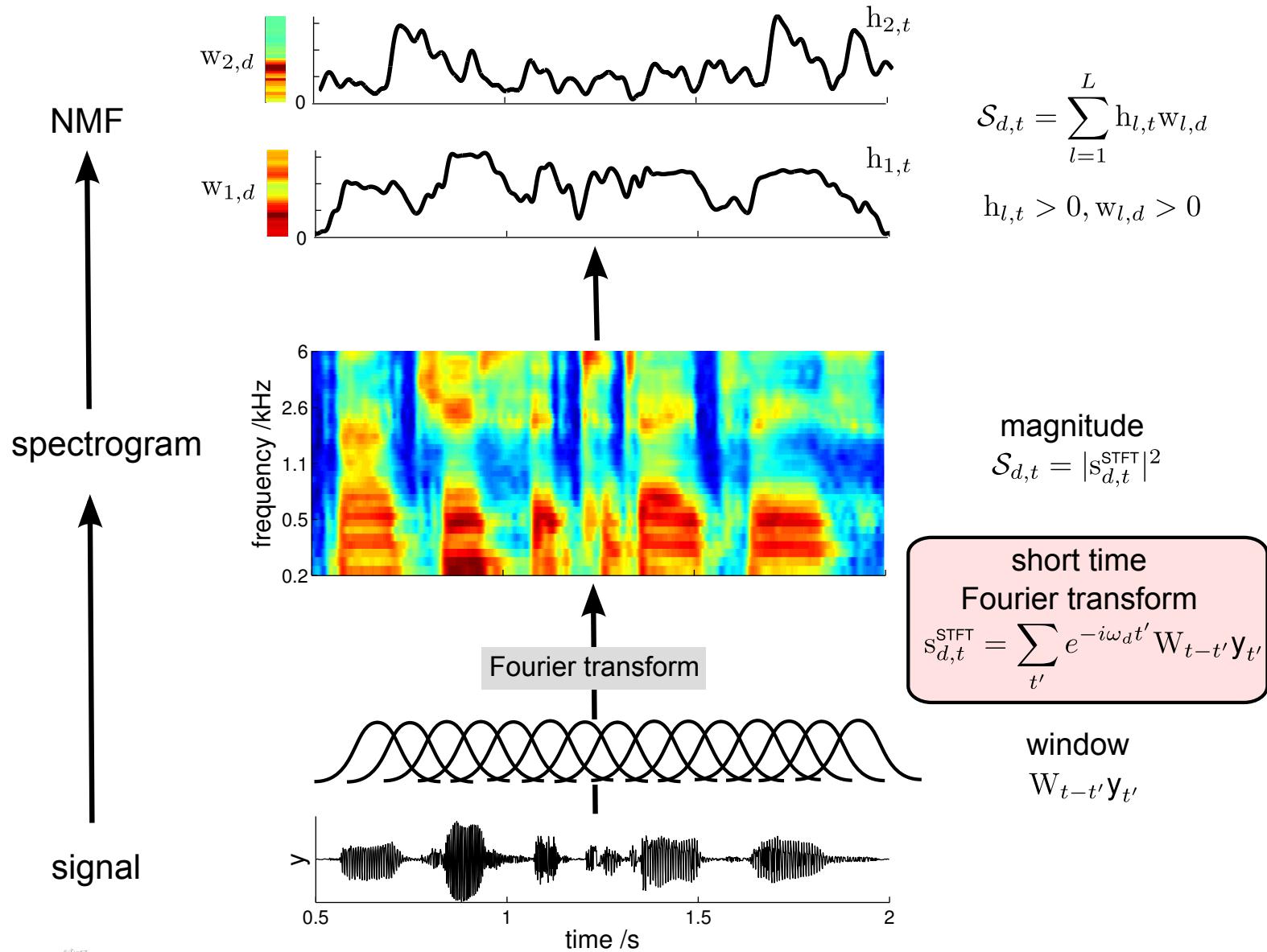
$$h_{l,t} > 0, w_{l,d} > 0$$

magnitude
 $s_{d,t} = |s_{d,t}^{\text{STFT}}|^2$

short time
Fourier transform
 $s_{d,t}^{\text{STFT}} = \sum_{t'} e^{-i\omega_d t'} W_{t-t'} y_{t'}$

window
 $W_{t-t'} y_{t'}$

A typical audio pipeline



What form of generative model corresponds to the STFT?

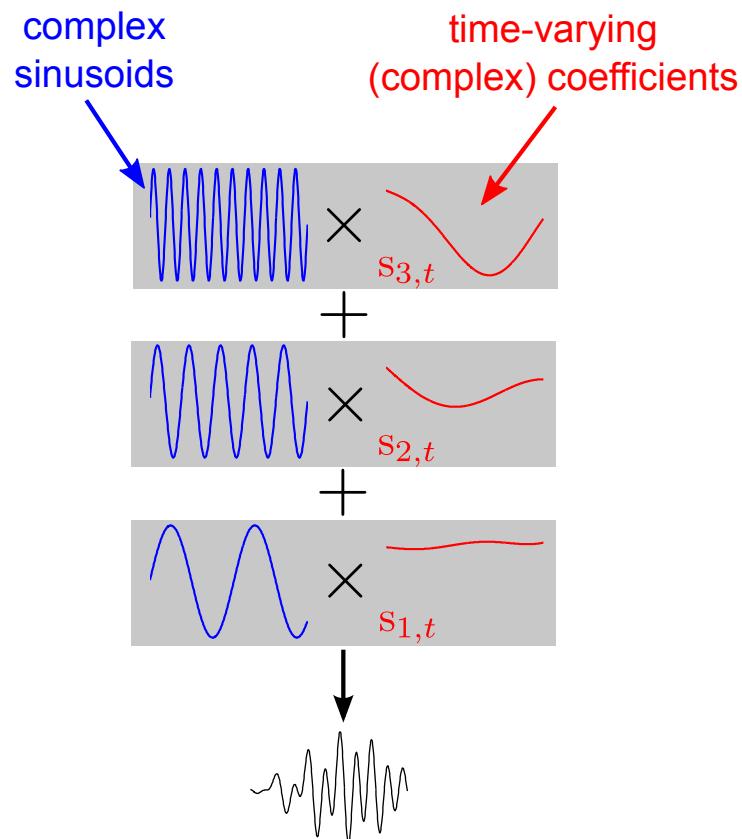
desire: expected value of latent time-frequency coefficients $s_{d,1:T} = \text{STFT}$

- assume y formed by (weighted) superposition of band-limited signals $s_{d,1:T}$
 - linearity of inference can be assured by setting the distributions of each $s_{d,1:T}$ and the noise to be Gaussian
 - time-invariance \implies generative model statistically stationary
- \implies GP prior over STFT coefficients, $p(s_{d,1:T}) = \mathcal{G}(s_{d,1:T}; 0, \Gamma)$, stationary

$$\Gamma_{t,t'} \approx \sum_{k=1}^T \text{FT}_{t,k}^{-1} \gamma_k \text{FT}_{k,t'} \text{ where } \text{FT}_{k,t} = e^{-2\pi i (k-1)(t-1)/T}$$

Time-frequency analysis as inference

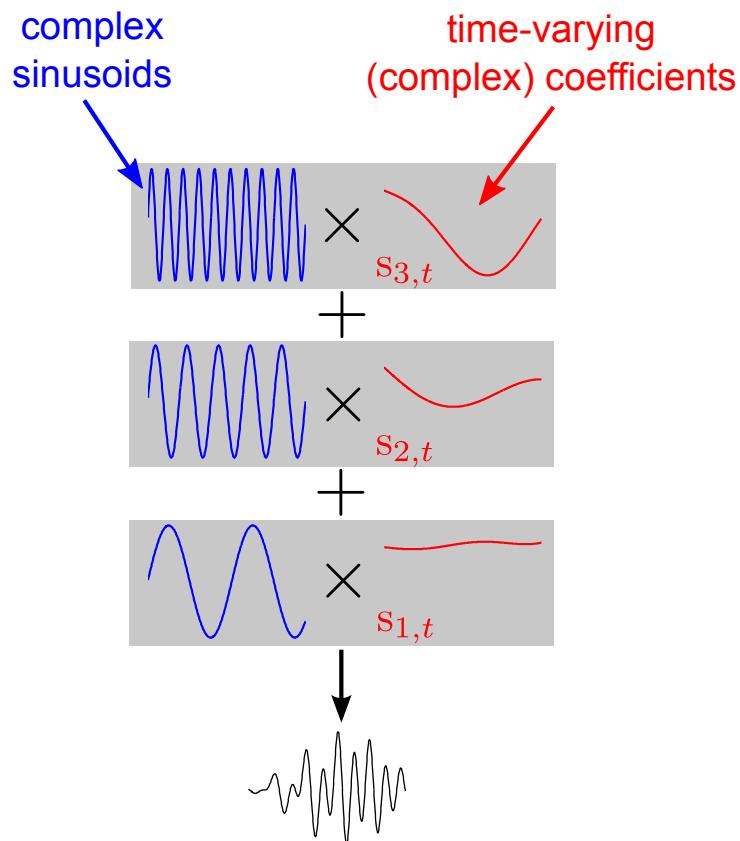
generation



$$y_t = \sum_d \Re(e^{i\omega_d t} s_{d,t}) + \sigma_y \eta_t$$

Time-frequency analysis as inference

generation

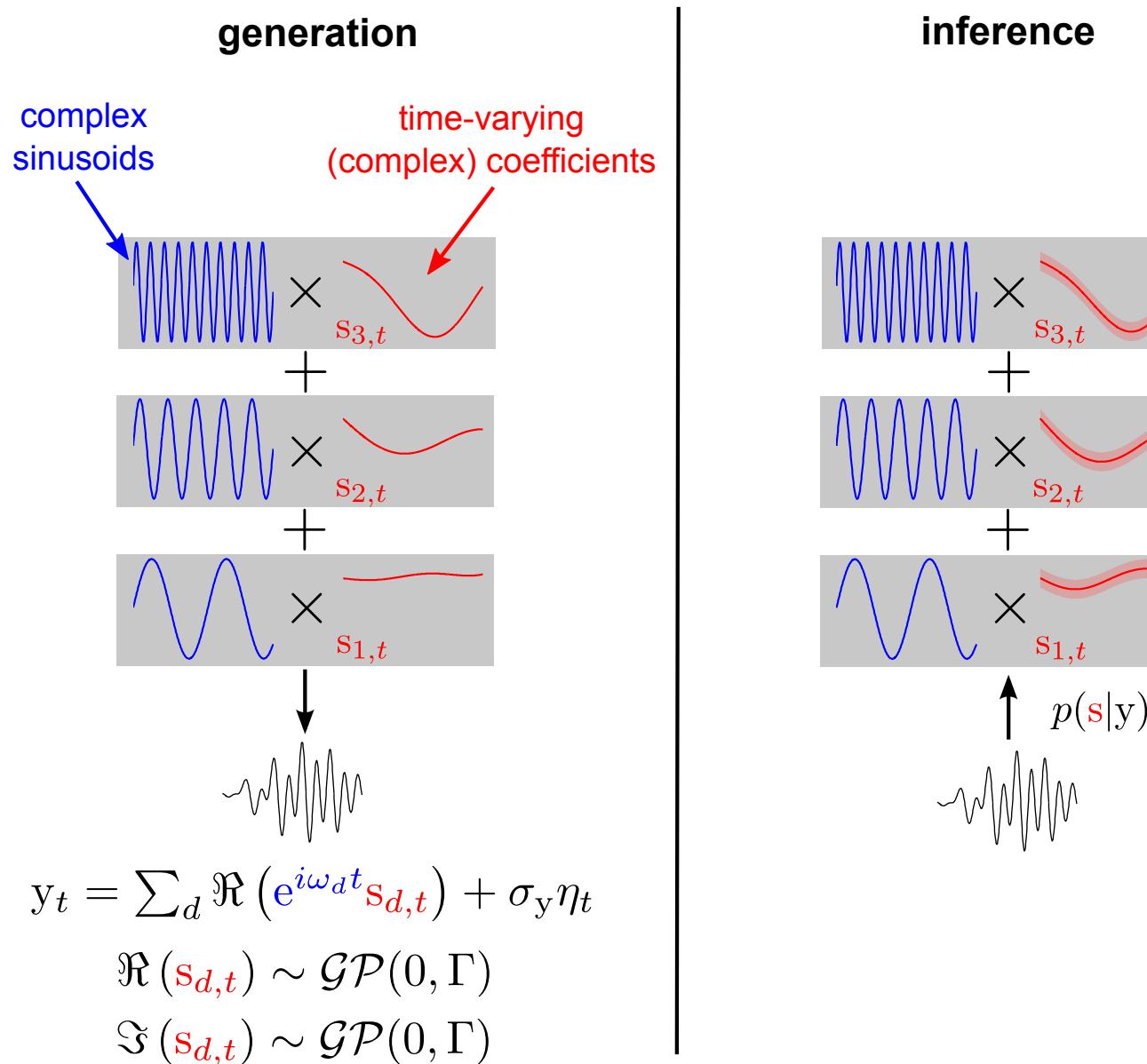


$$y_t = \sum_d \Re(e^{i\omega_d t} s_{d,t}) + \sigma_y \eta_t$$

$$\Re(s_{d,t}) \sim \mathcal{GP}(0, \Gamma)$$

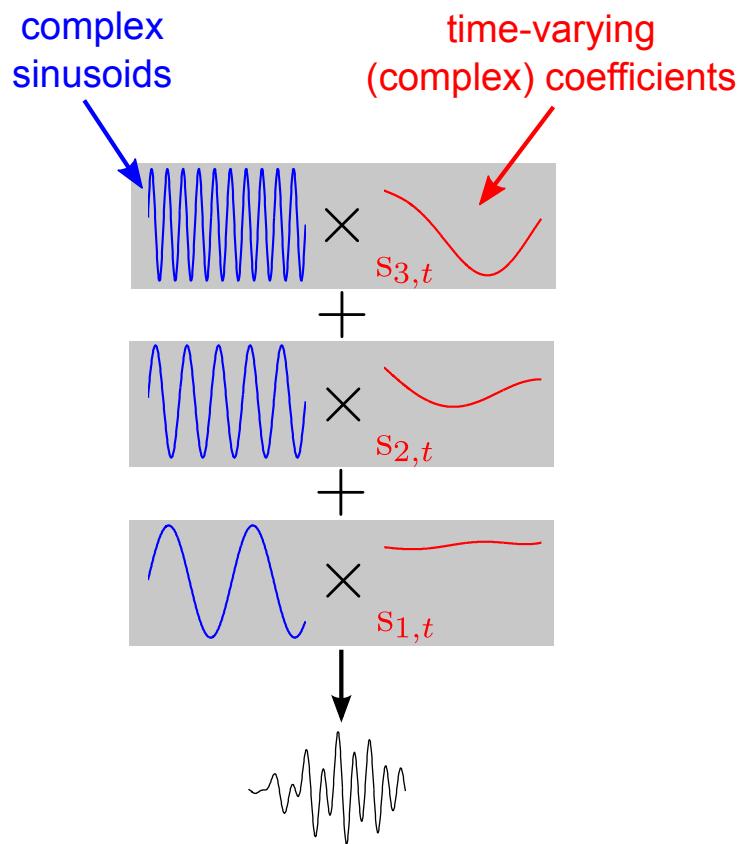
$$\Im(s_{d,t}) \sim \mathcal{GP}(0, \Gamma)$$

Time-frequency analysis as inference



Time-frequency analysis as inference

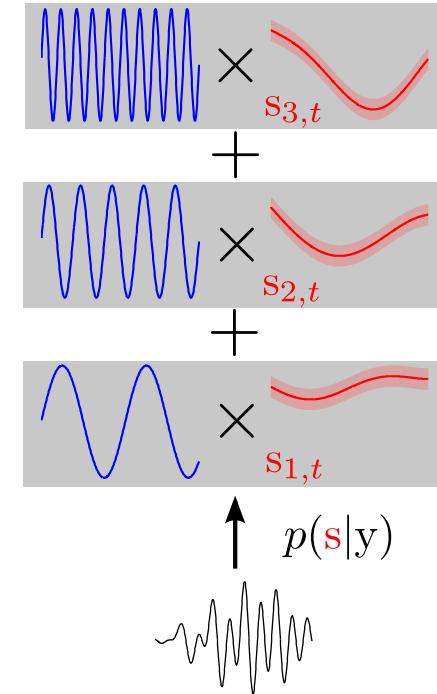
generation



inference

most probable coefficients given the signal is the STFT

$$\text{STFT} \equiv \arg \max_s p(\mathbf{s}|y)$$



STFT window = prior covariance

\otimes
frequency shifted
inverse signal covariance

Time-frequency analysis as inference

generation

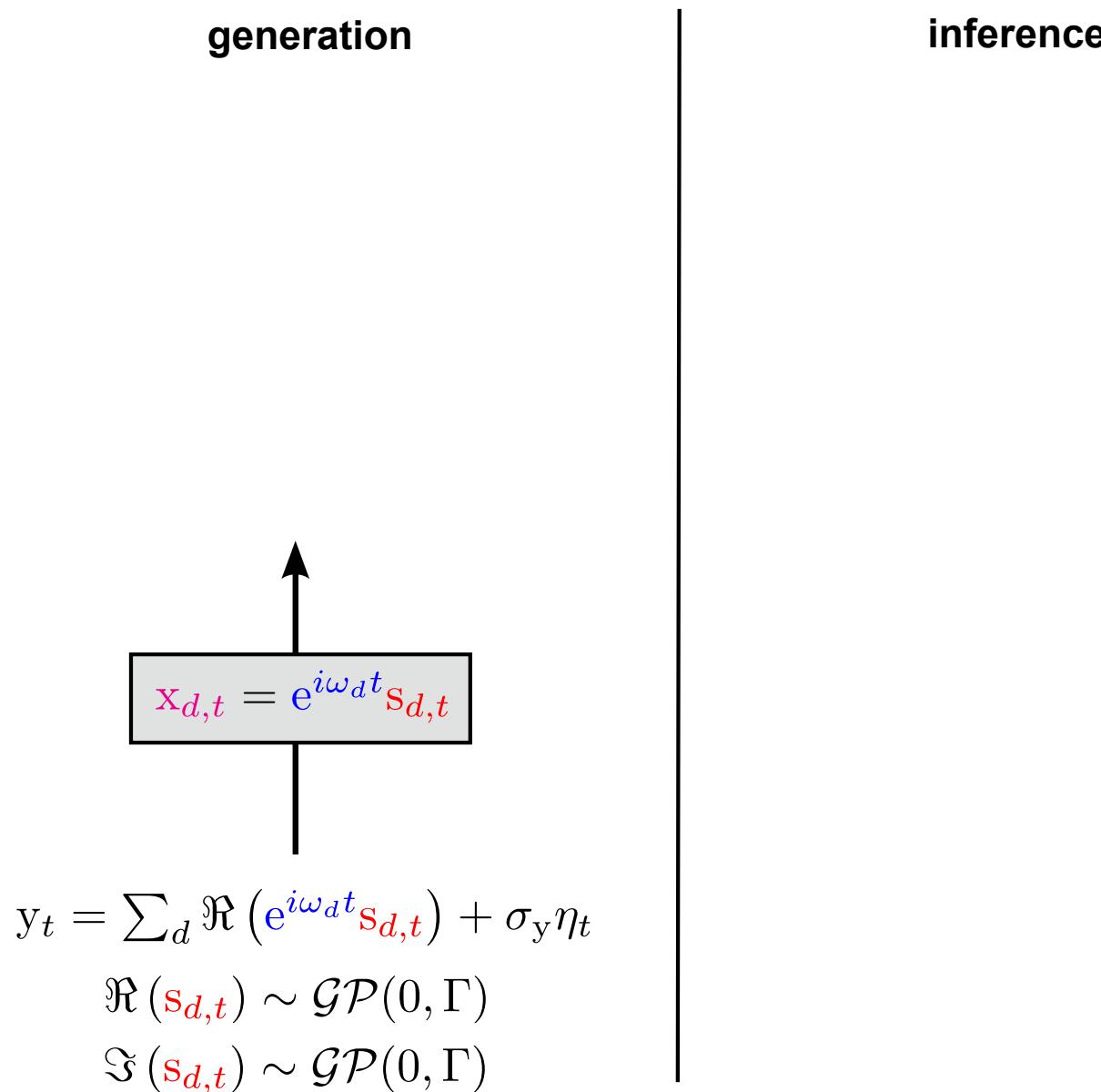
$$y_t = \sum_d \Re(e^{i\omega_d t} s_{d,t}) + \sigma_y \eta_t$$

$$\Re(s_{d,t}) \sim \mathcal{GP}(0, \Gamma)$$

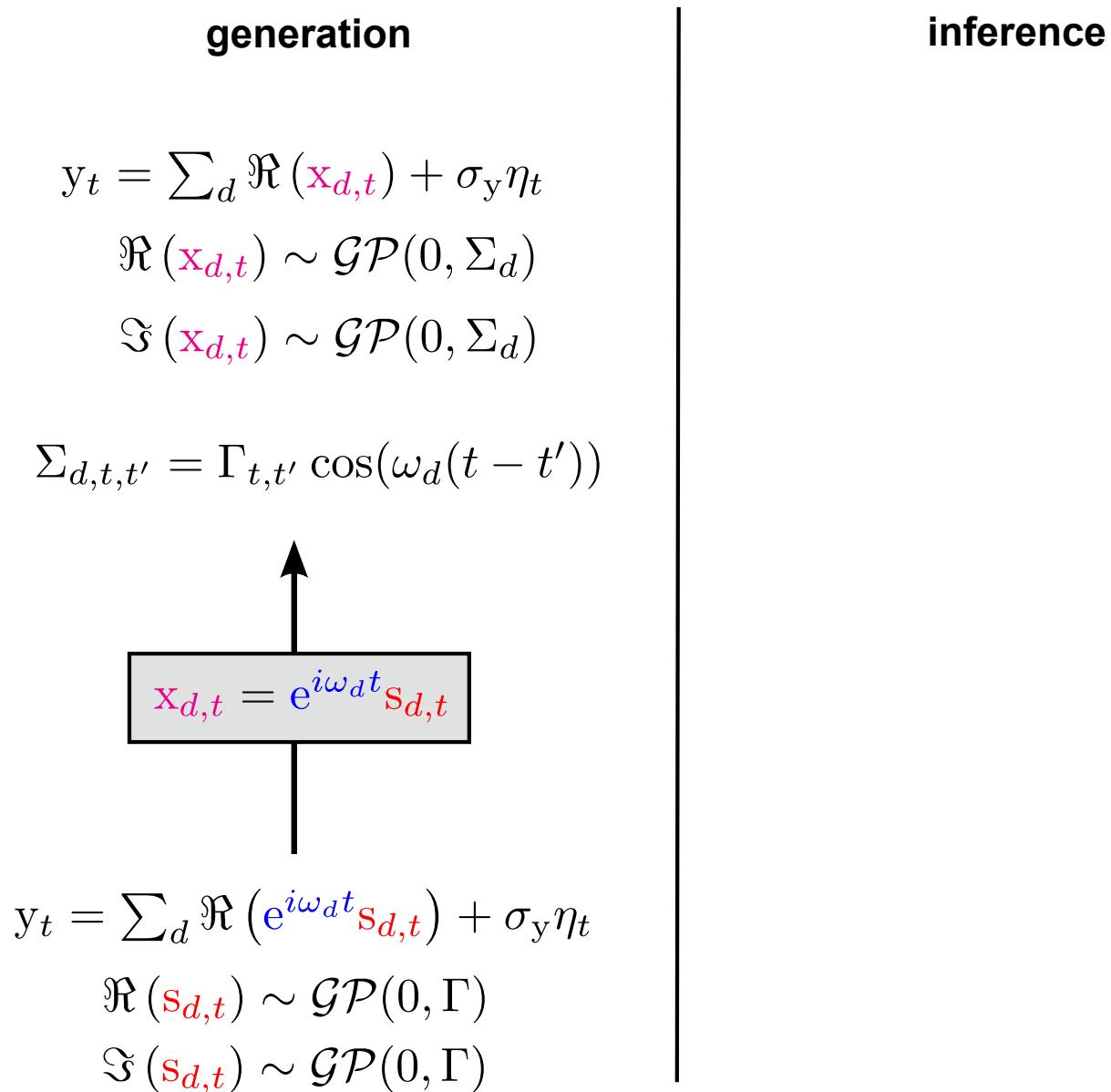
$$\Im(s_{d,t}) \sim \mathcal{GP}(0, \Gamma)$$

inference

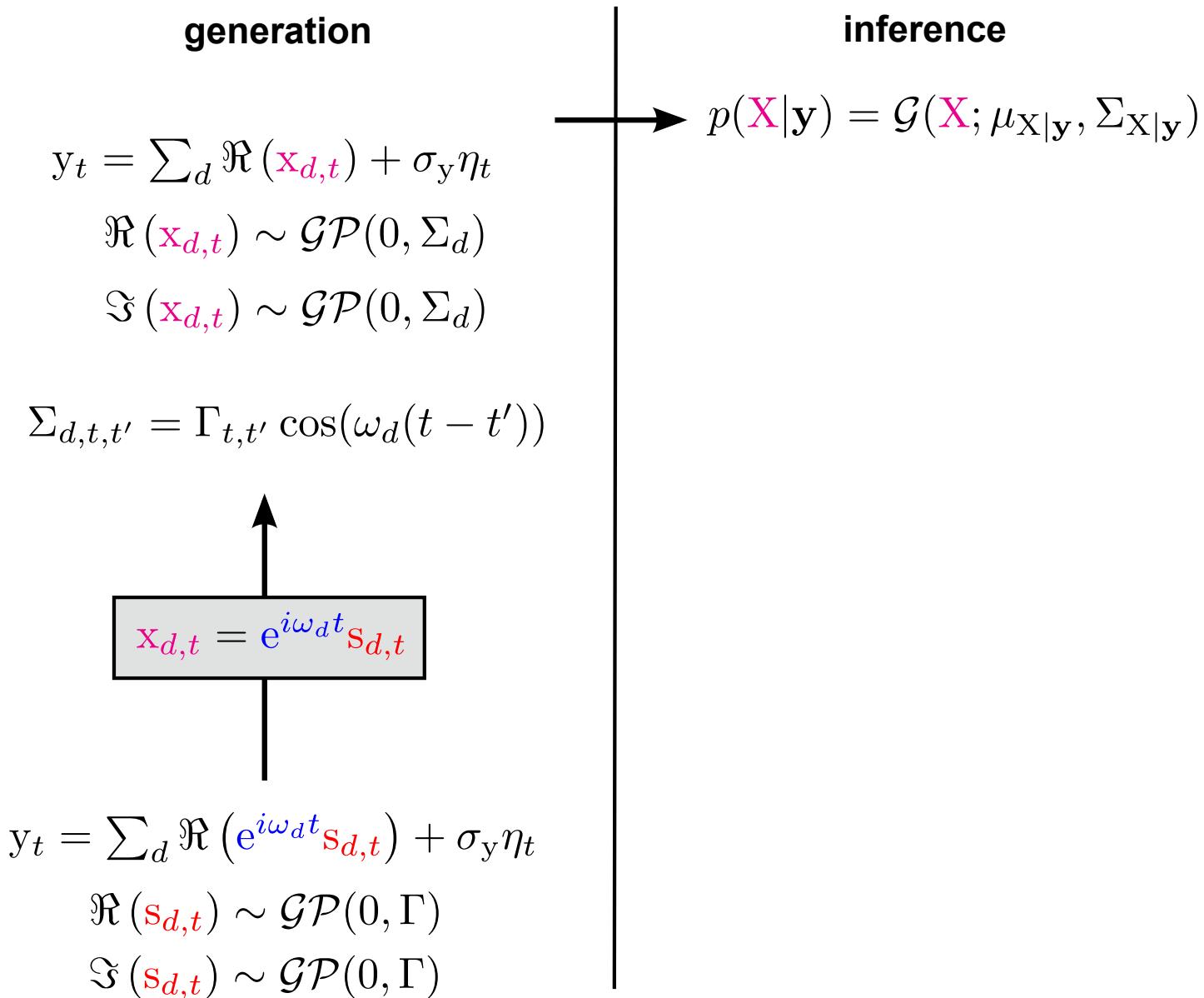
Time-frequency analysis as inference



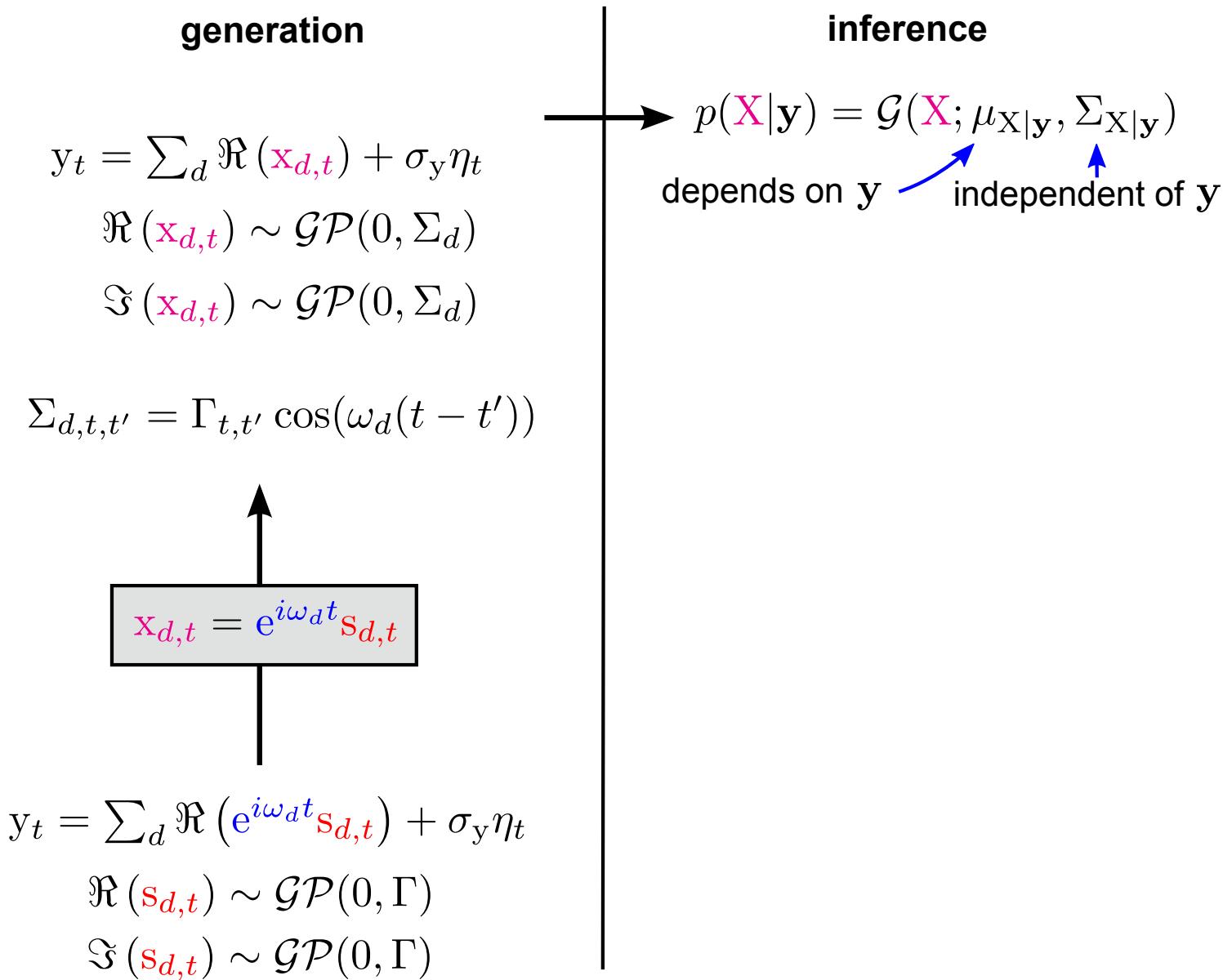
Time-frequency analysis as inference



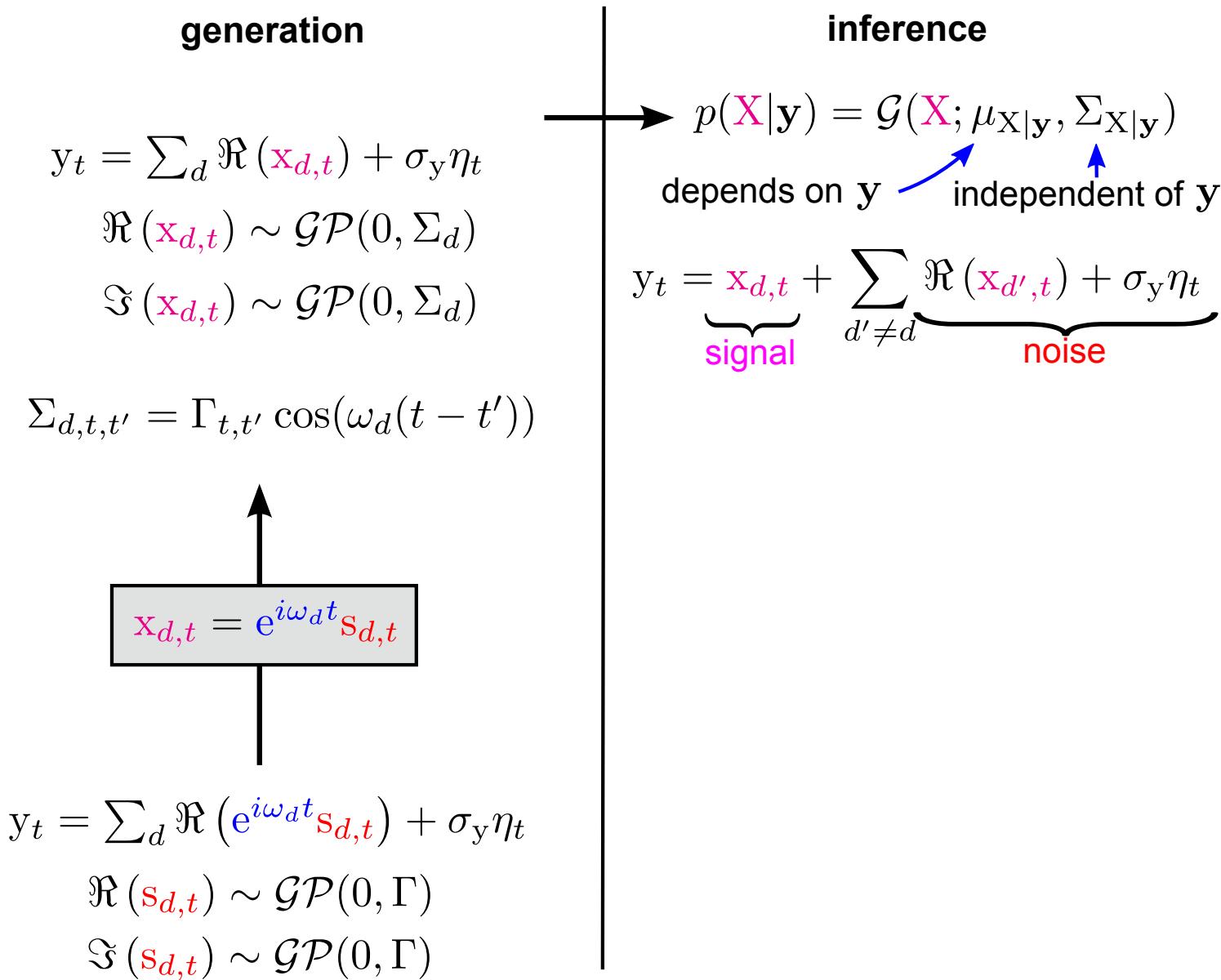
Time-frequency analysis as inference



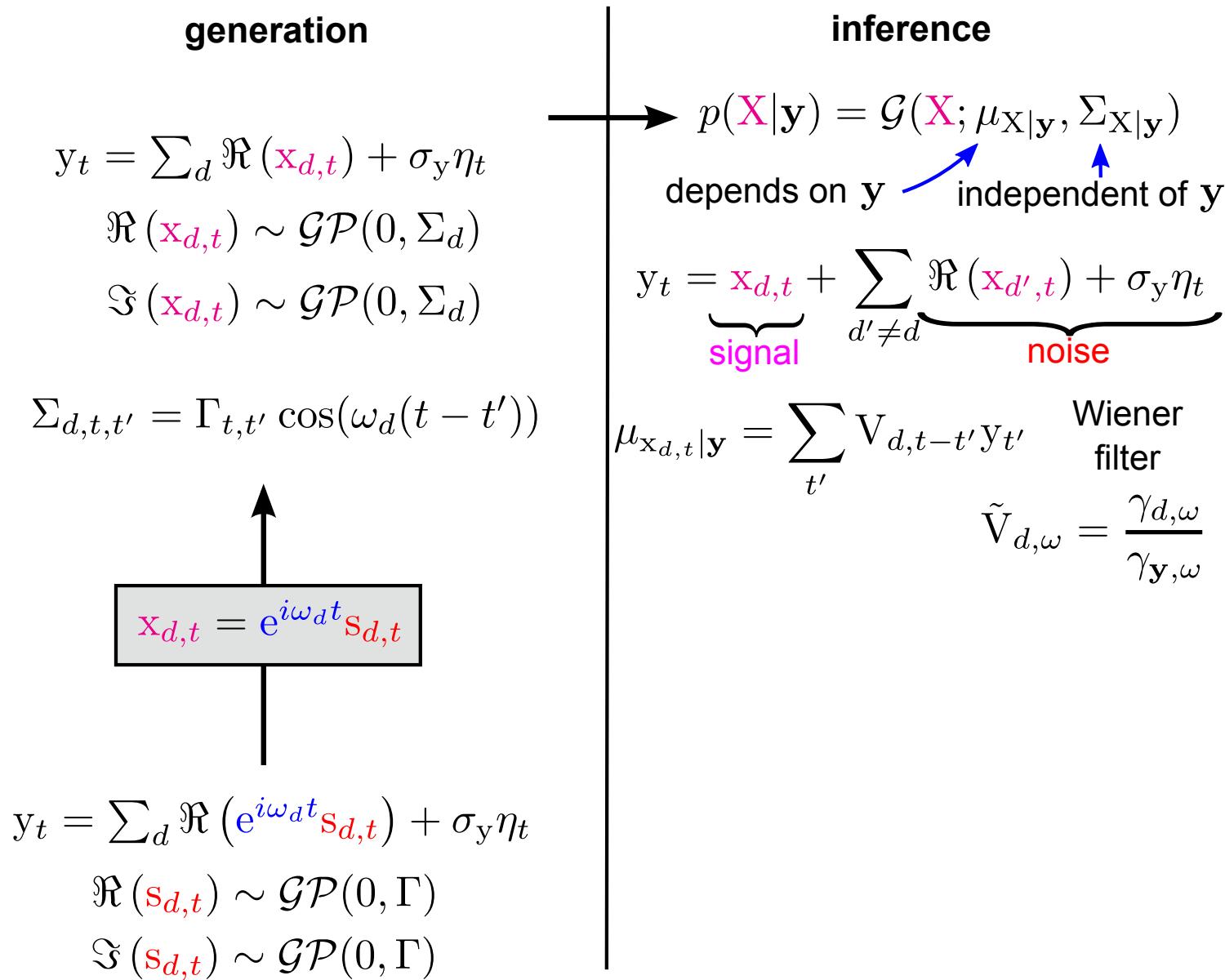
Time-frequency analysis as inference



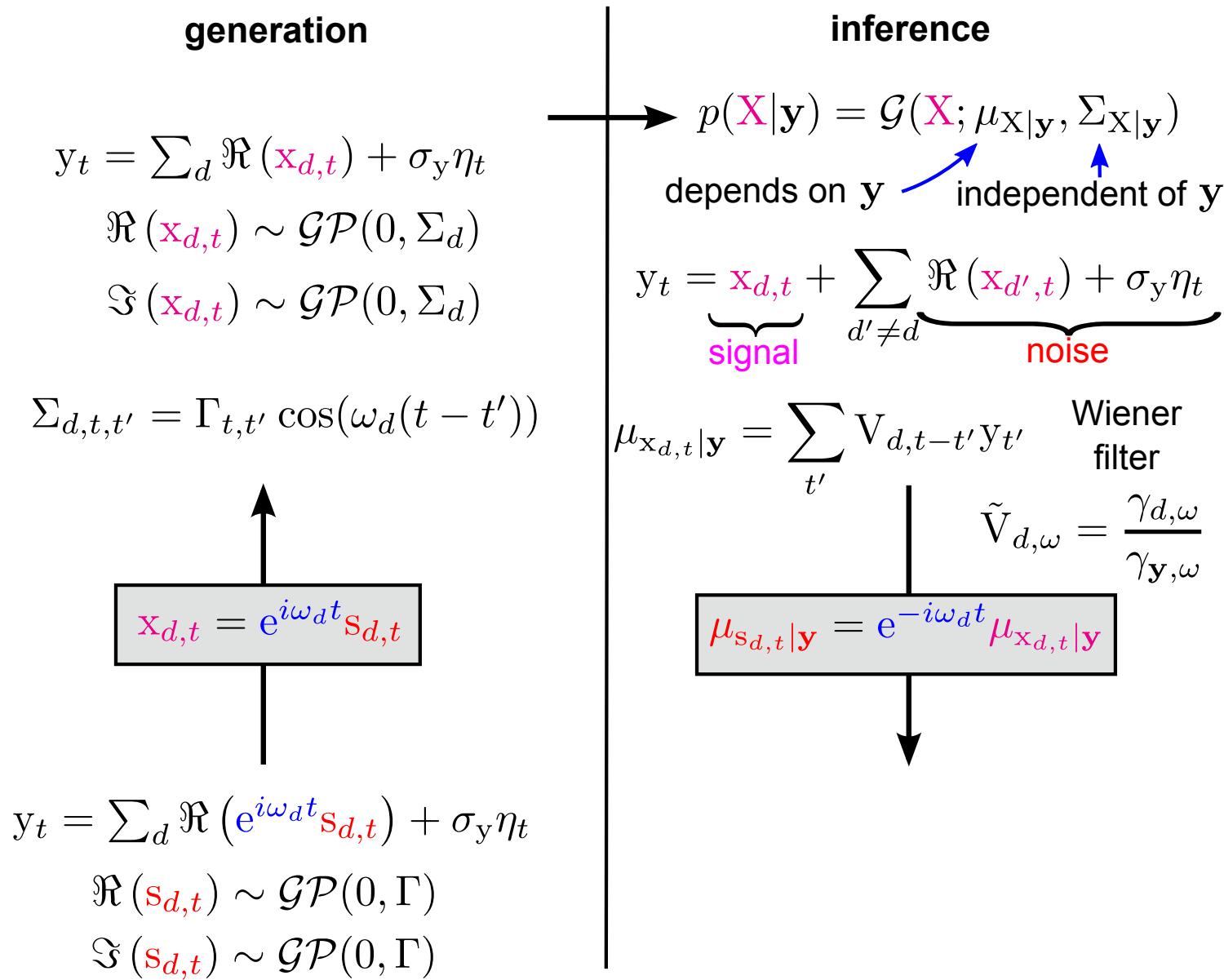
Time-frequency analysis as inference



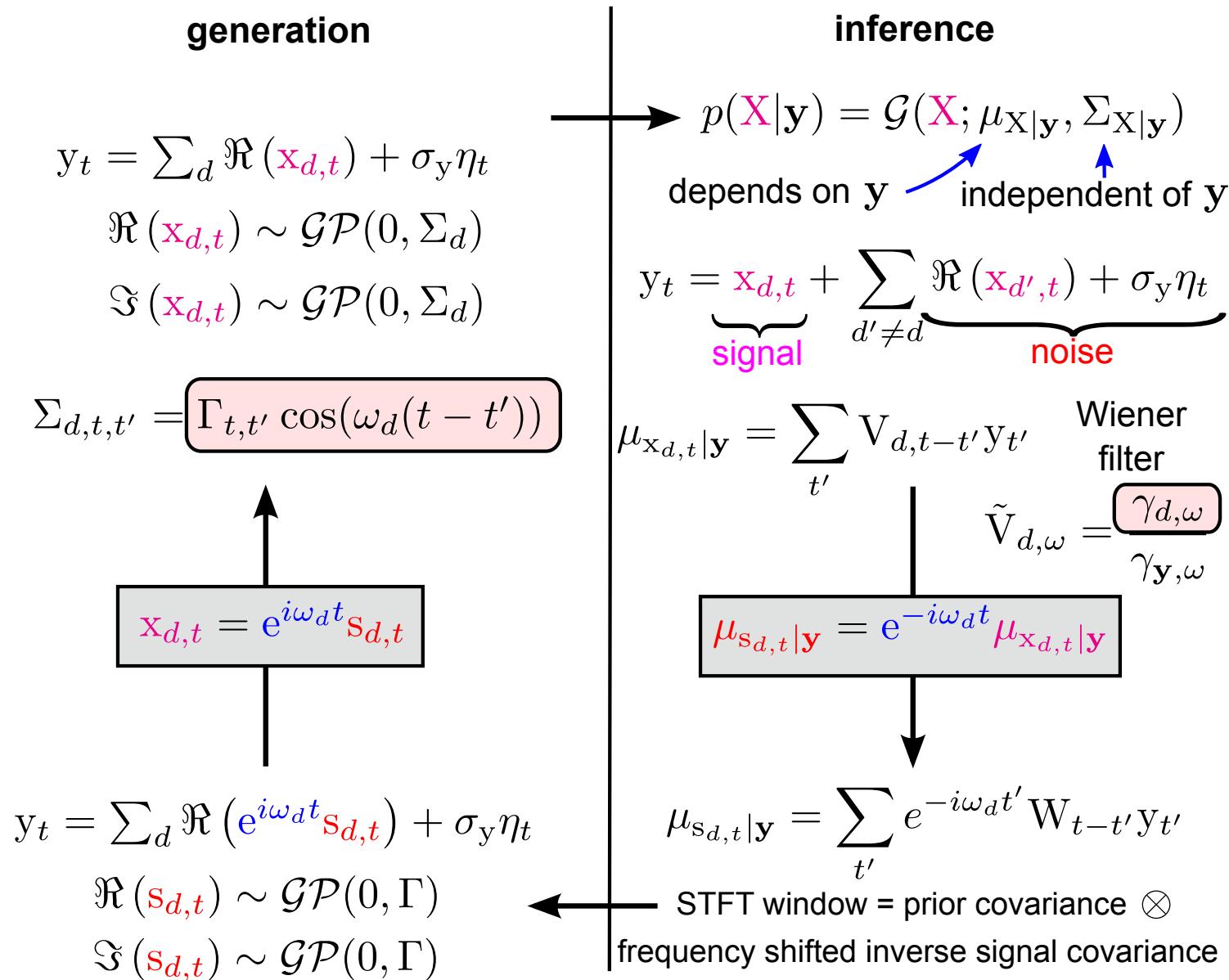
Time-frequency analysis as inference



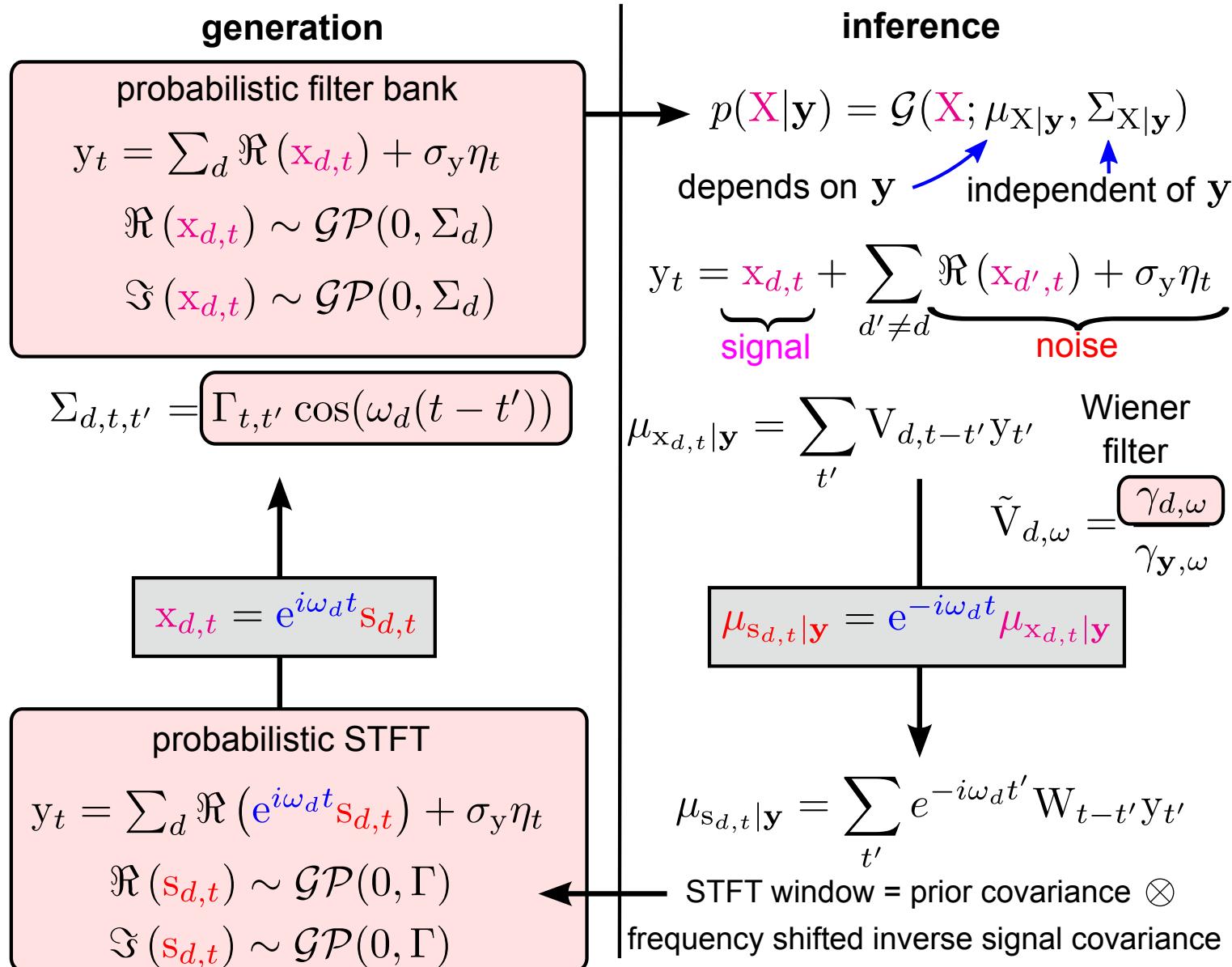
Time-frequency analysis as inference



Time-frequency analysis as inference



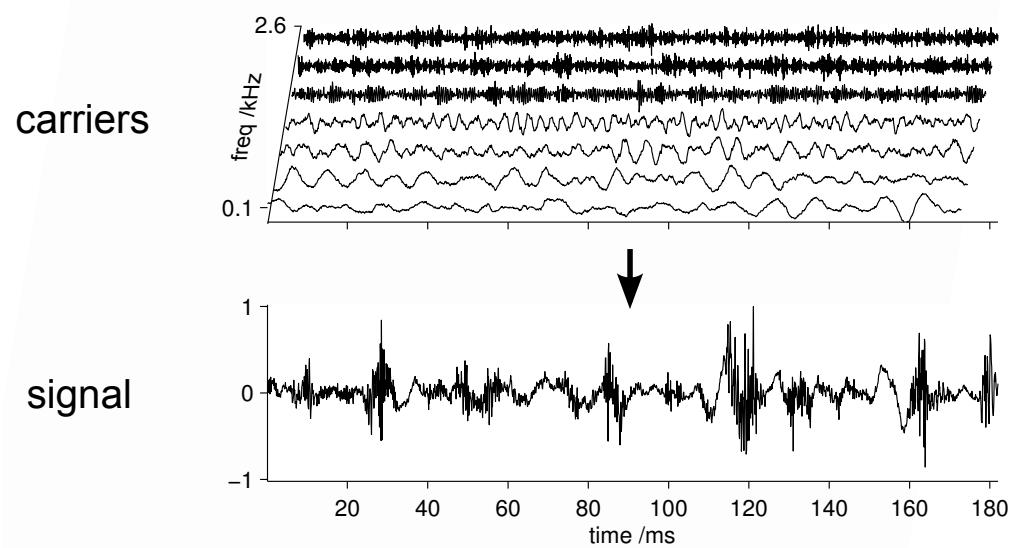
Time-frequency analysis as inference



Time-frequency analysis as inference

- probabilistic models in which **inference recovers STFT, filter bank, wavelet analysis**
 - **unifies a number of existing probabilistic time-series models** & connects to traditional sig. proc.
 - can learn window of STFT and frequencies (equivalently filter properties)
 - frequency shift relationship mimics classical relationship between these time-frequency relationships
- **hops/down-sampling and finite window used correspond to FITC** (uniformly spaced pseudo-points) and sparse-covariance approximations
 - rediscover Nyquist in the context of approximation GPs

Probabilistic audio processing pipeline

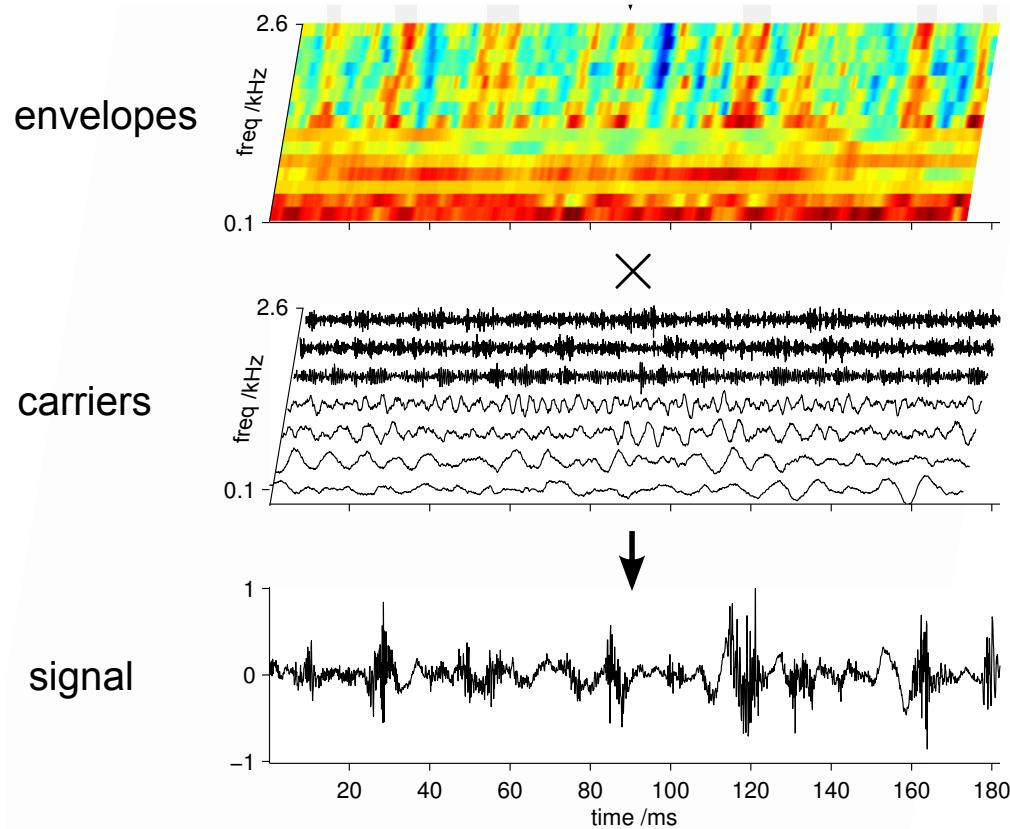


$$c_d(t) \sim \text{GP} \left(0, \frac{\text{mean}}{0} \frac{\text{spectrum}}{f} \right)$$

= bandpass
Gaussian
noise

$$y(t) = \sum_{d=1}^D \Re(x_d(t))$$

Probabilistic audio processing pipeline



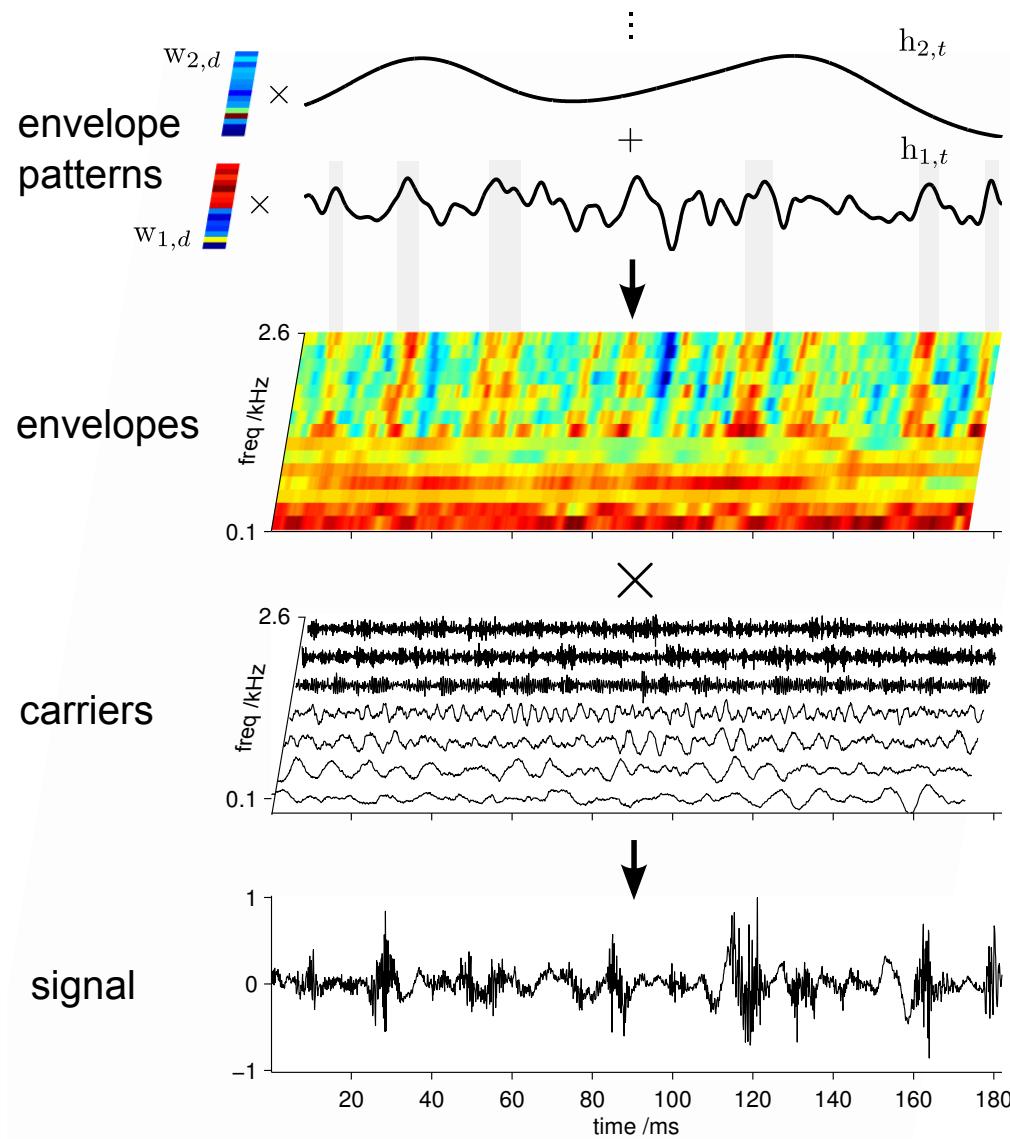
$a_{d,t}$

$$c_d(t) \sim GP\left(0, \frac{\text{mean}}{f} \text{ spectrum} \right)$$

= bandpass
Gaussian
noise

$$y(t) = \sum_{d=1}^D \Re(x_d(t)) a_{d,t}$$

Probabilistic audio processing pipeline



mean spectrum
 $\log h_{l,t} \sim \text{GP}\left(\mu_k, \frac{f}{0} f\right)$
= slow Gaussian process

$$a_{d,t} = \sum_{l=1}^L h_{l,t} w_{l,d}$$

mean spectrum
 $c_d(t) \sim \text{GP}\left(0, \frac{f}{0} f\right)$
= bandpass Gaussian noise

$$y(t) = \sum_{d=1}^D \Re(x_d(t)) a_d(t)$$

Inference and Learning

- **Key Observation** – fix envelopes:
 - posterior over carriers is Gaussian
 - posterior mean given by an (adaptive) filter
- Leads to MAP estimation of the envelopes (or HMCMC), let $z_{lt} = \log h_{lt}$

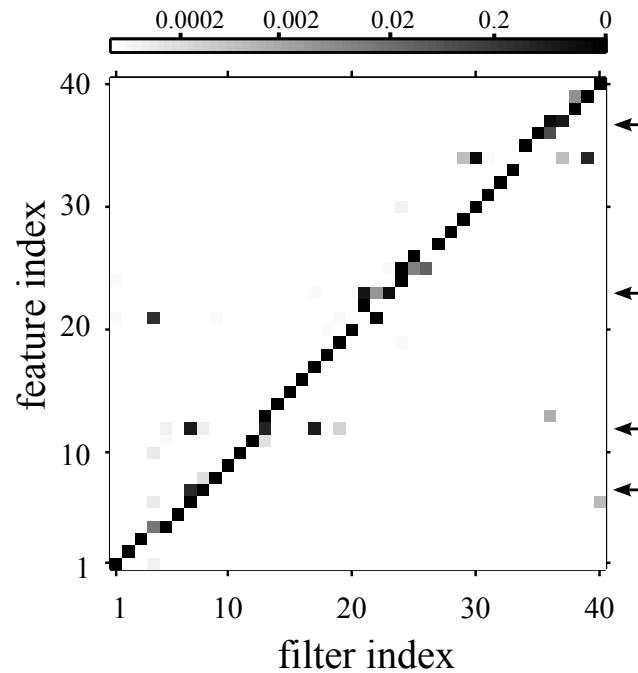
$$Z^{\text{MAP}} = \arg \max_Z p(Z|Y)$$

$$p(Z|Y) = \frac{1}{Z} p(Z, Y) = \frac{1}{Z} \int d\mathbf{X} p(Z, Y, \mathbf{X}) = \frac{1}{Z} p(Z) \int d\mathbf{X} p(Y|\mathbf{A}, \mathbf{X}) p(\mathbf{X})$$

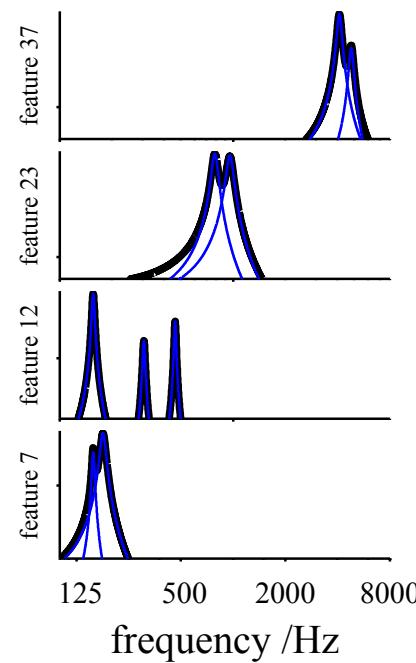
- Compute integral efficiently using **chain stuctured approximation** and **Kalman Smoothing**
- Leads to gradient based optimisation for transformed amplitudes
- Learning: approximate Maximum Likelihood $\theta = \arg \max_{\theta} p(Y|\theta)$
- **NMF**: zero-temperature EM, one E-Step, initialise constant envelopes

Audio modelling

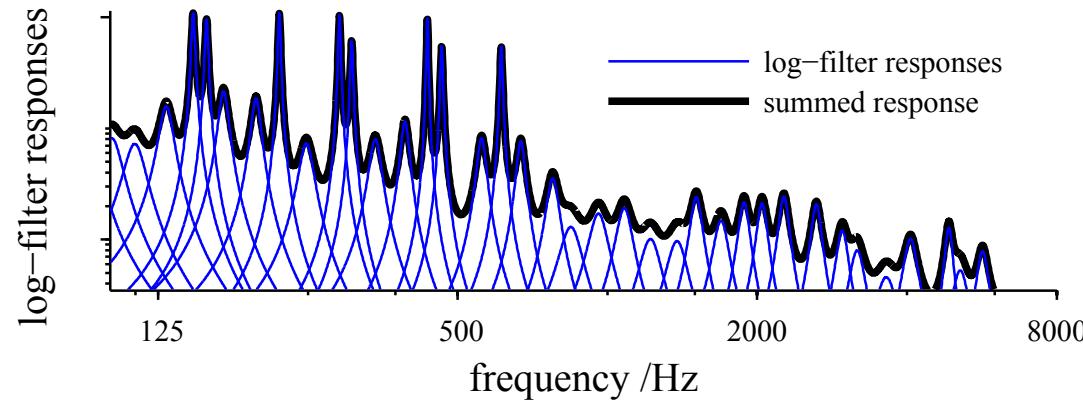
A. spectral features W



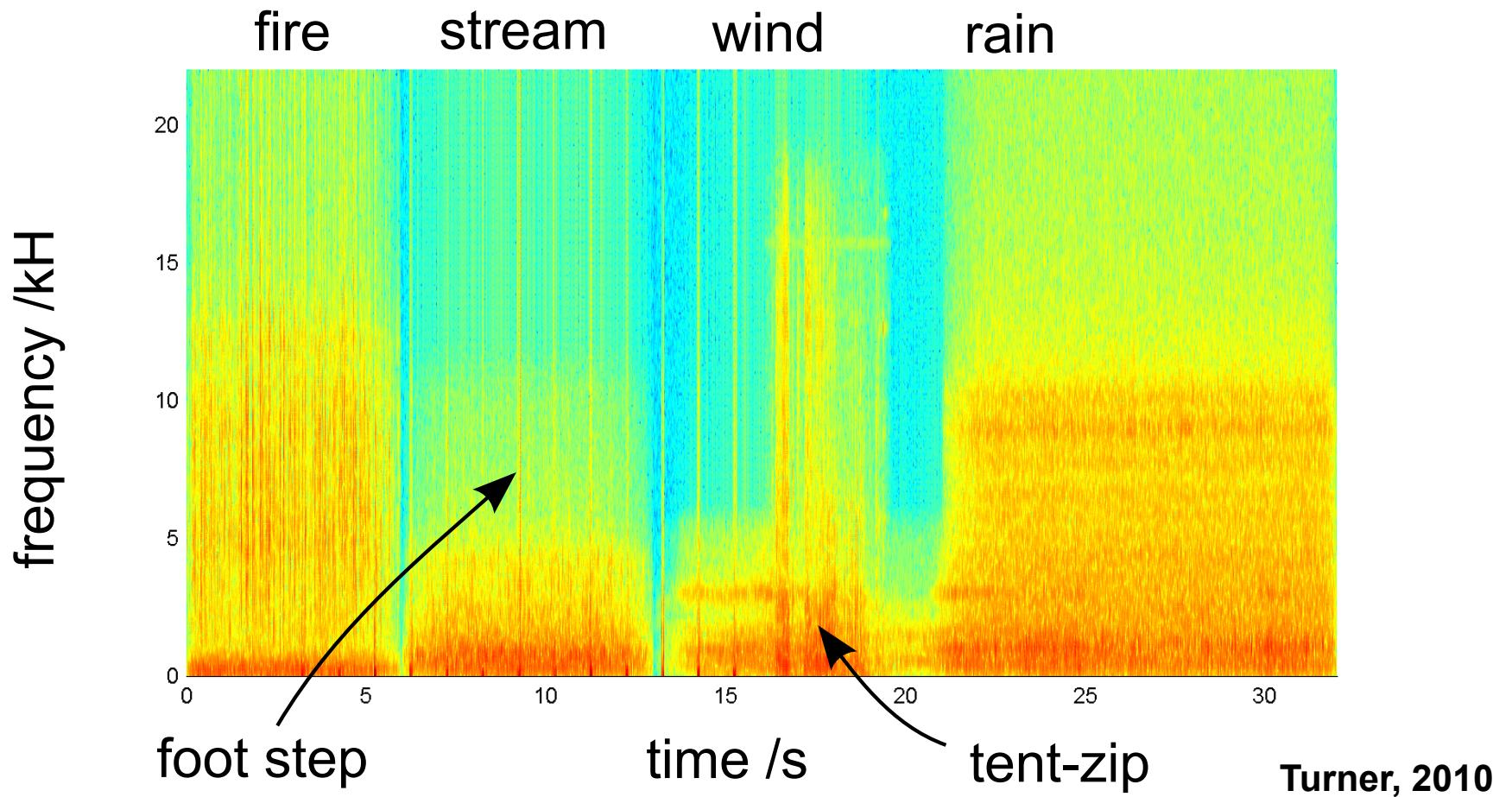
C. example features



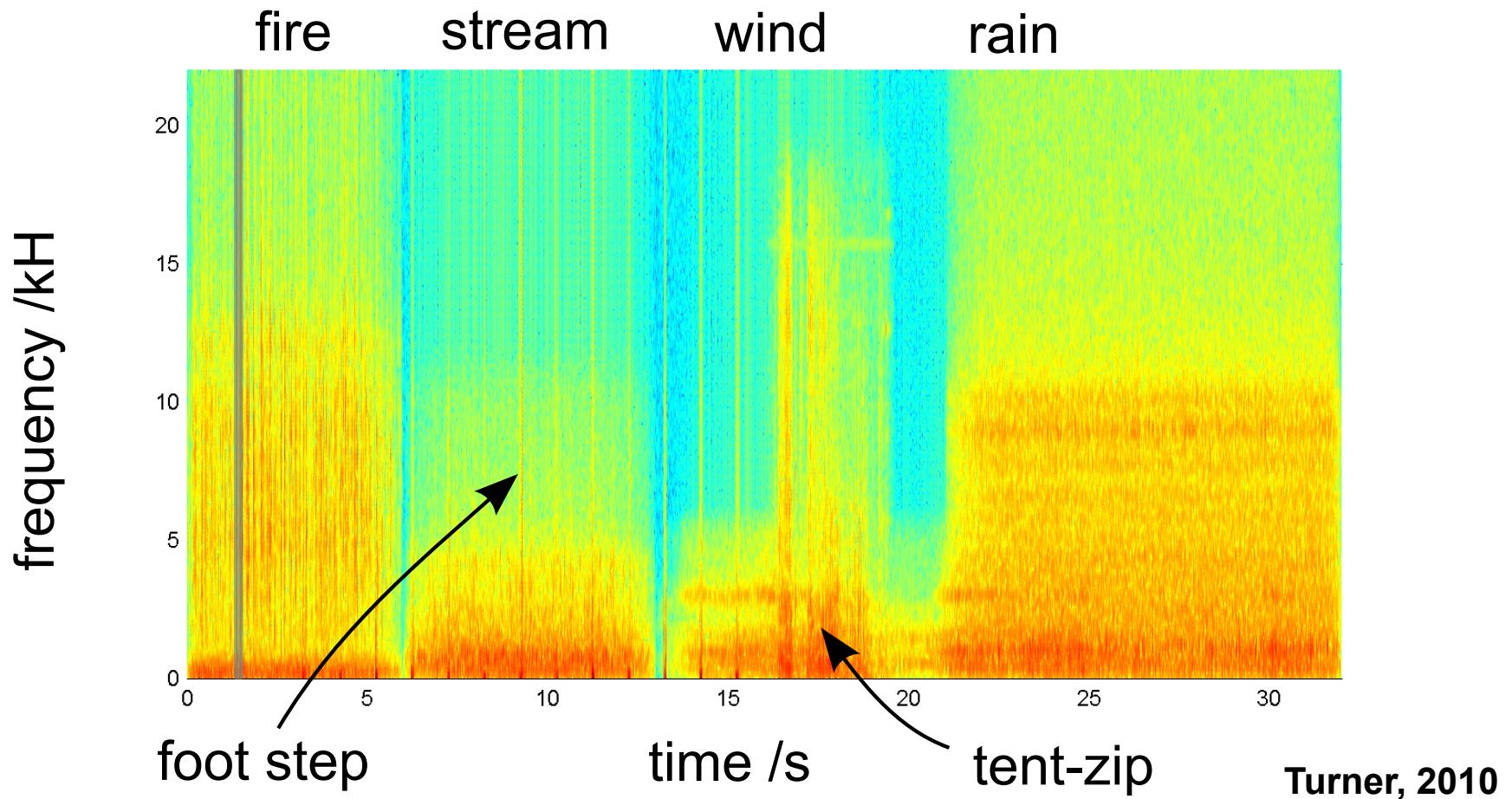
B. filterbank $\gamma(\theta)$



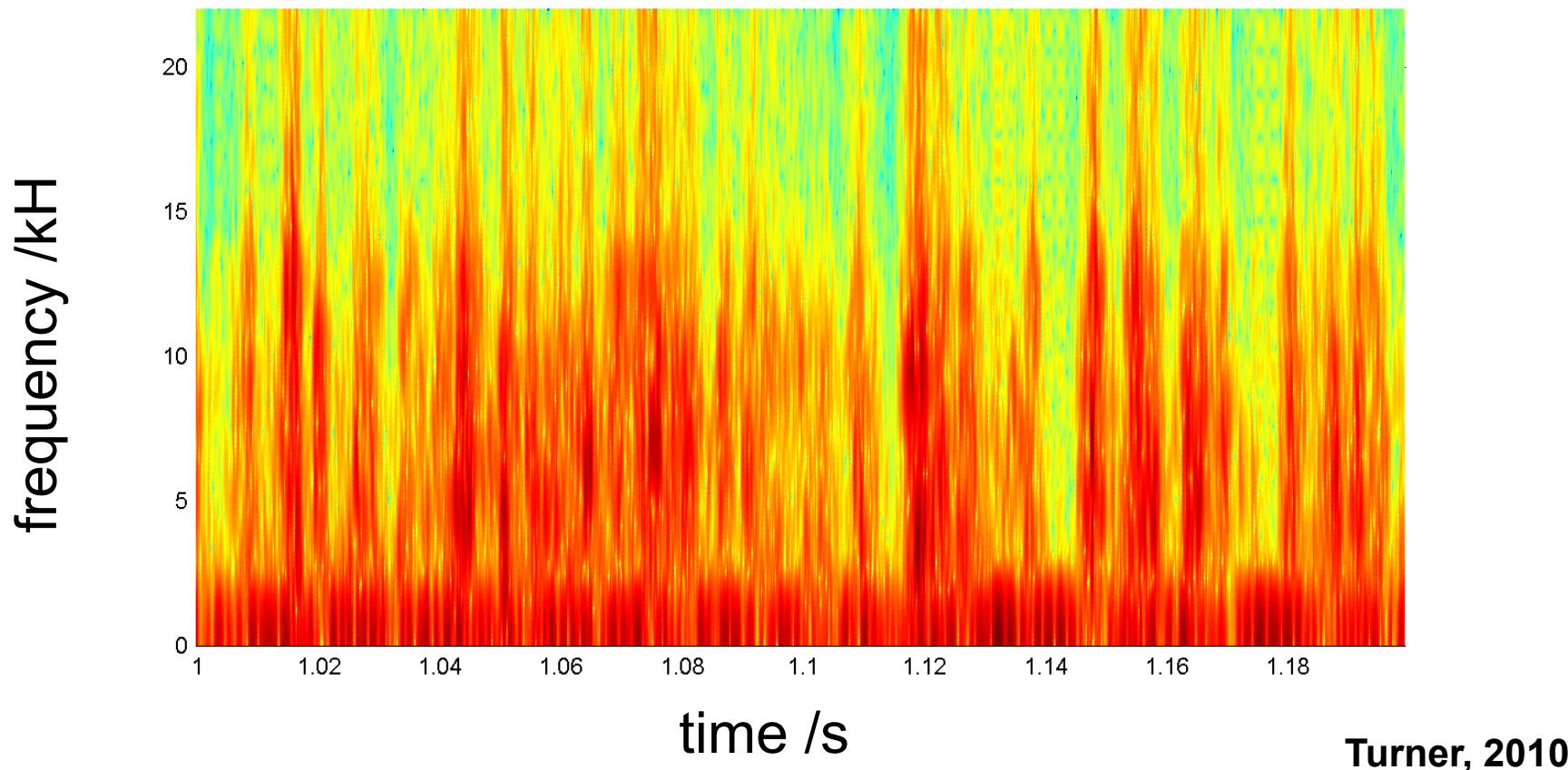
Audio modelling



Audio modelling



Audio modelling

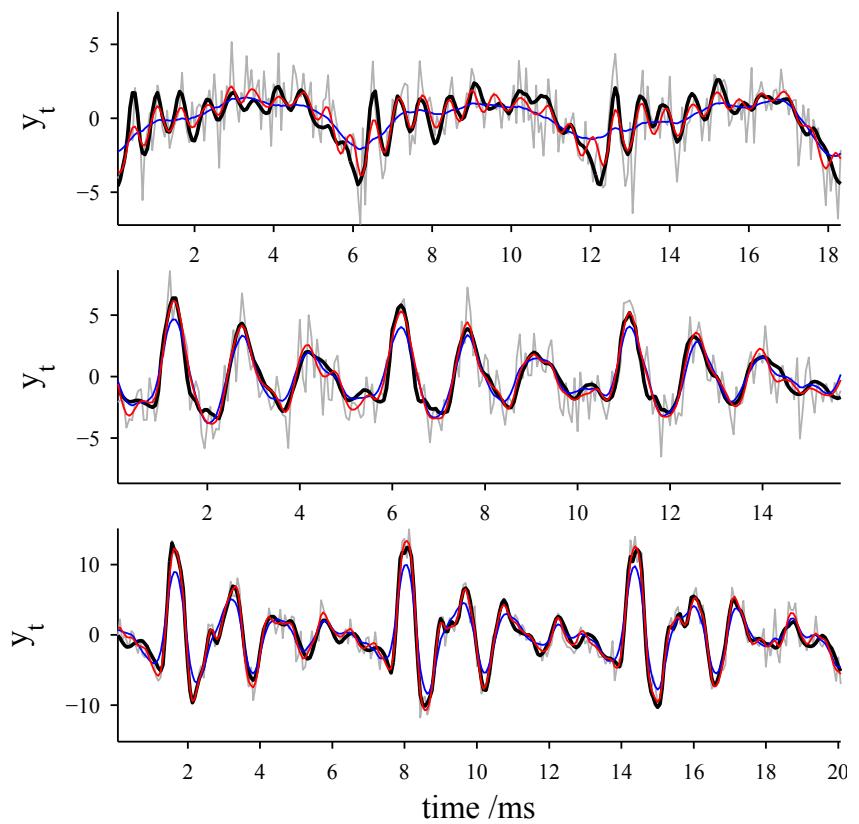
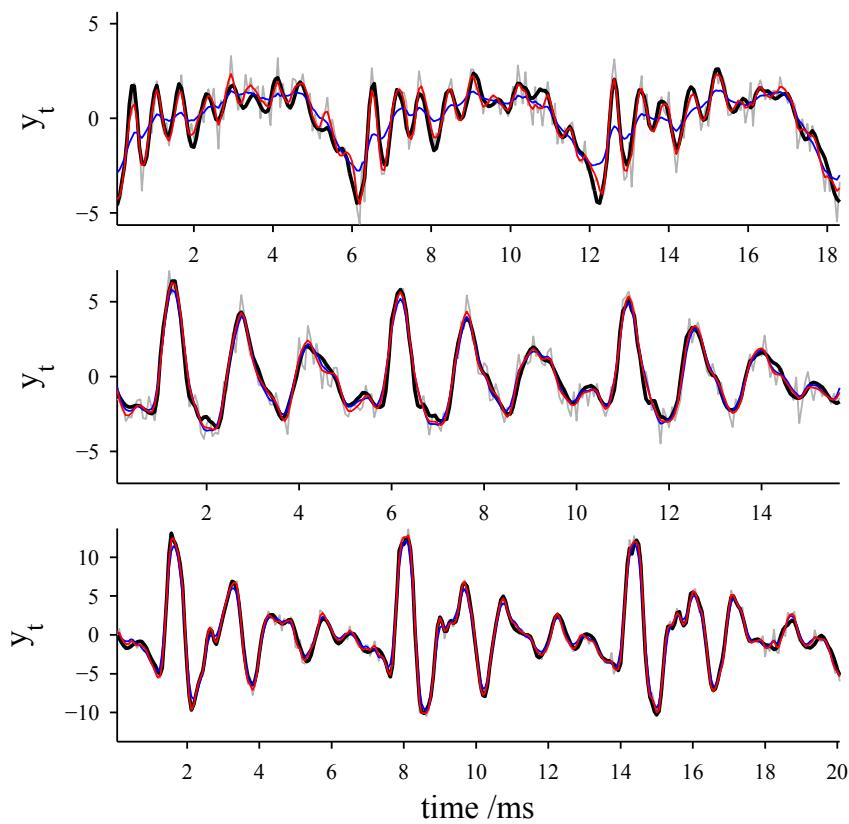
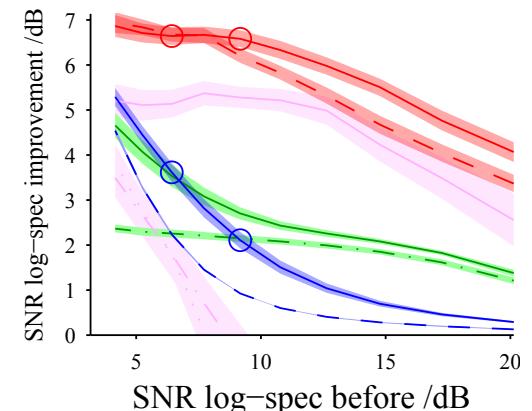
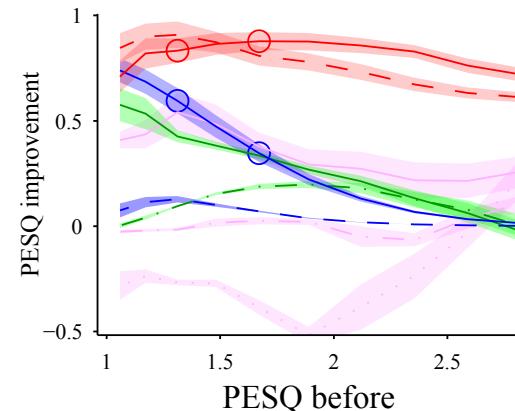
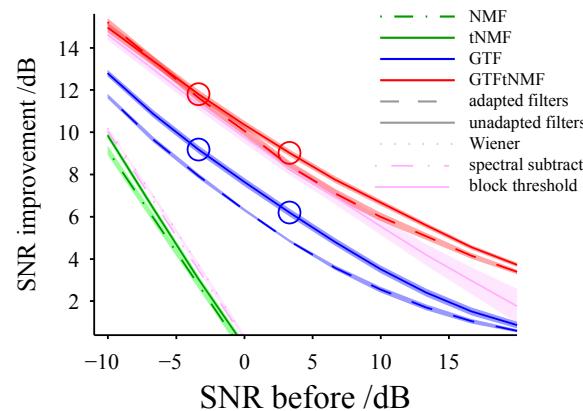


Turner, 2010

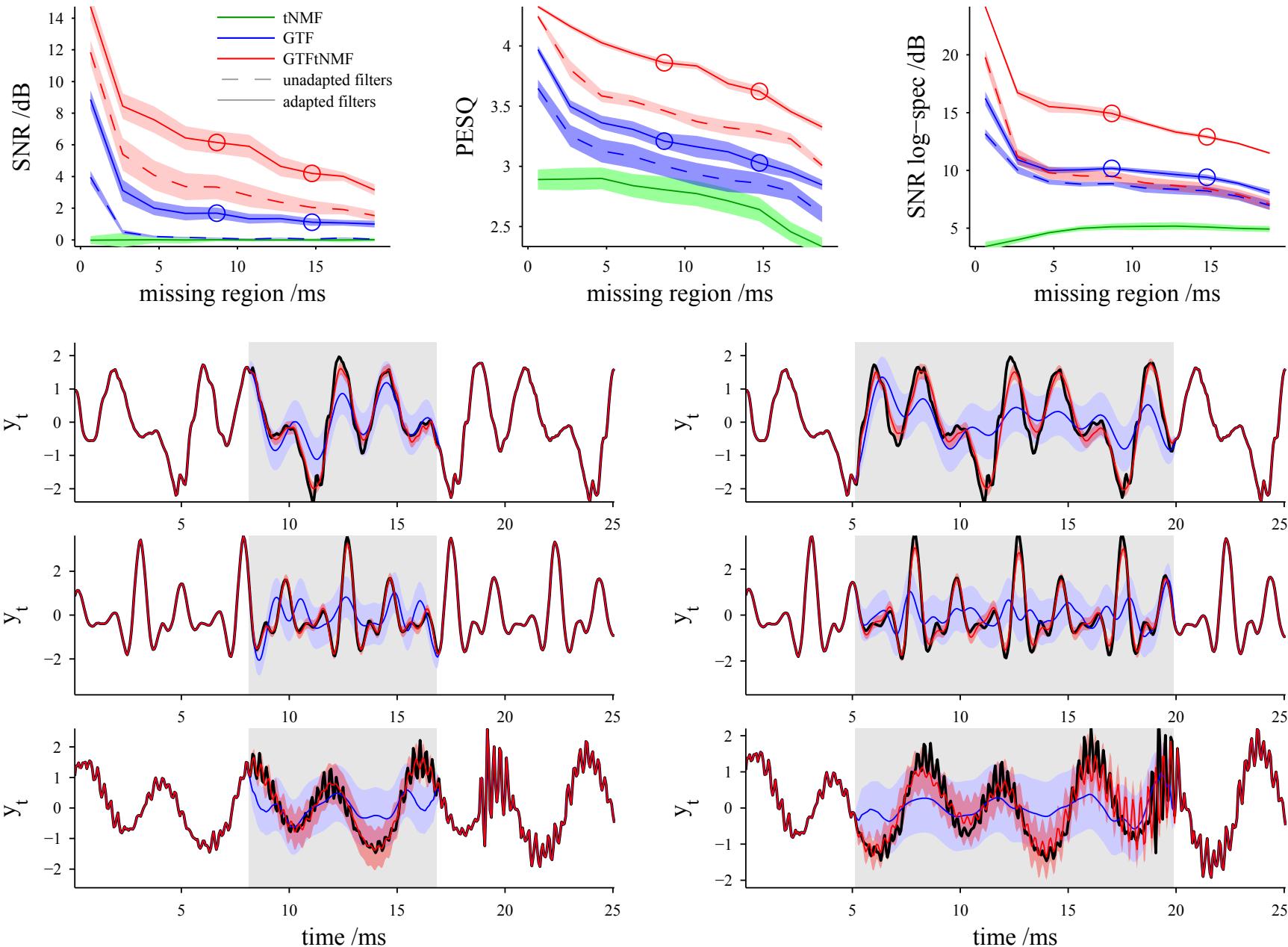
Statistical texture synthesis

- Old approach: build **detailed physical models** (e.g. rain drops)
- New approach
 - **train model** on your favourite texture
 - **sample** from the prior, and then from the likelihood.
- Waveform unique, but statistically matched to original
- Often perceptually indistinguishable

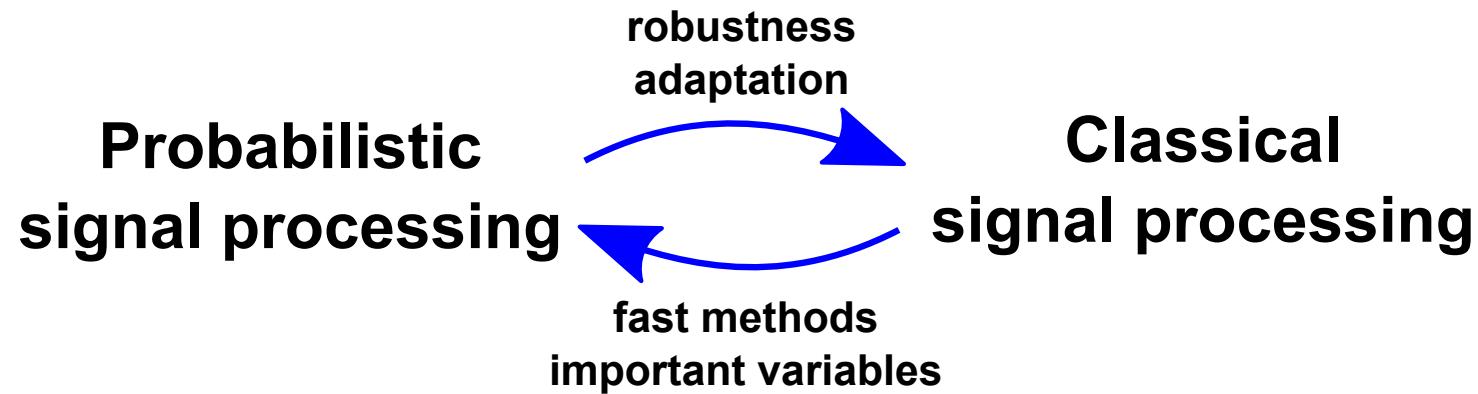
Audio denoising



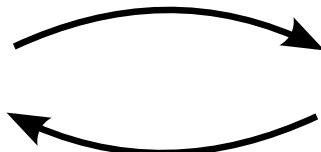
Audio missing data imputation



Unifying classical and probabilistic audio signal processing

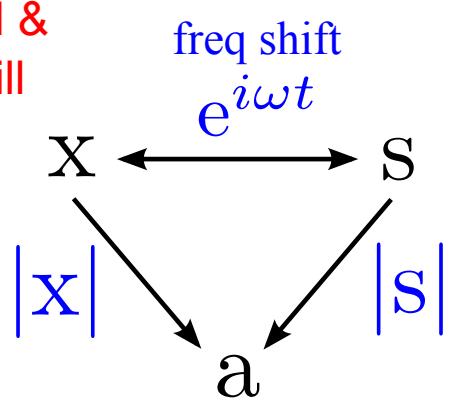


Probabilistic signal processing



Classical signal processing

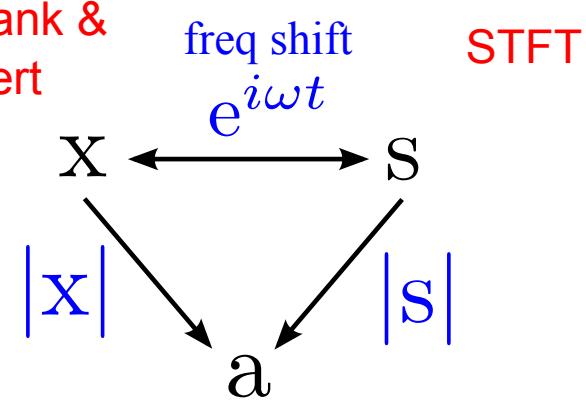
Cemgil &
Godsill



Qi &
Minka

equiv
estimation

Filter Bank &
Hilbert



Amplitudes

Spectrogram

Additional slides