

*ECE209AS (Winter 2021)*

# Lecture 5: Learning with Irregularly Sampled Time Series Data

---

Mani Srivastava

mbs@ucla.edu

Networked & Embedded Systems Lab

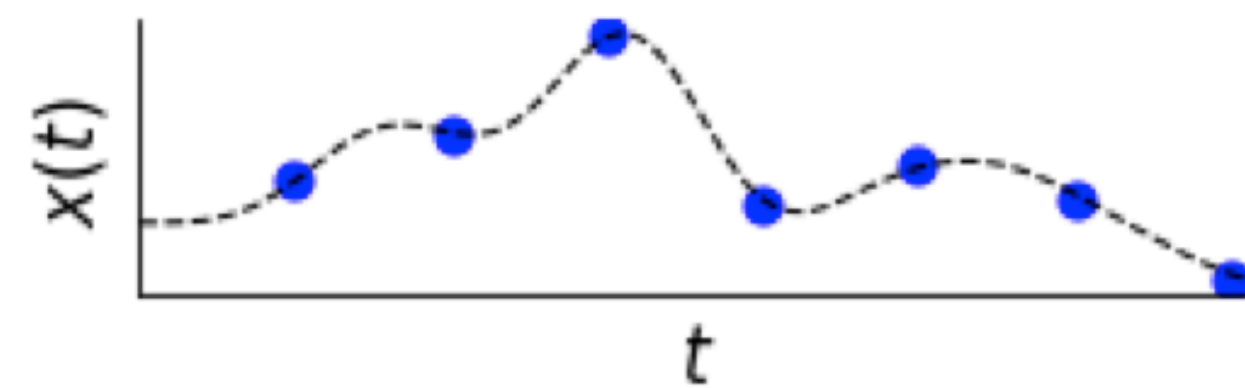
ECE & CS Departments

UCLA

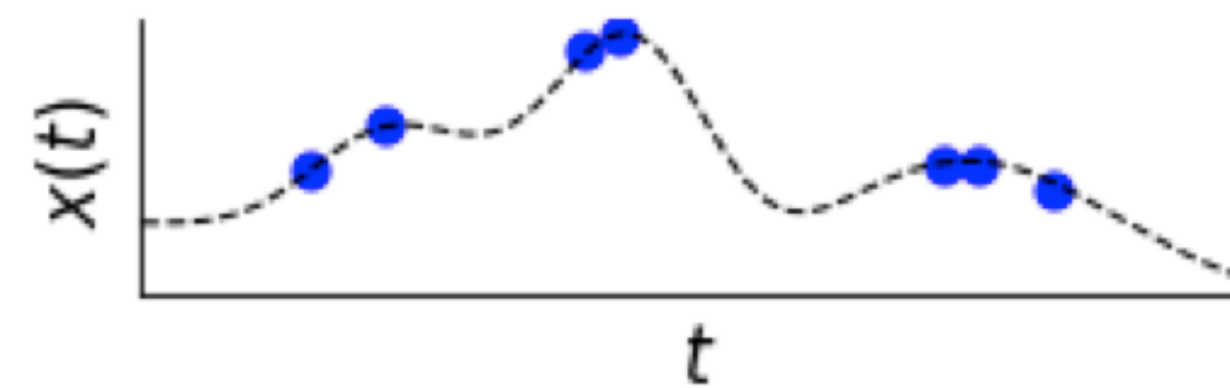


**Samueli**  
School of Engineering

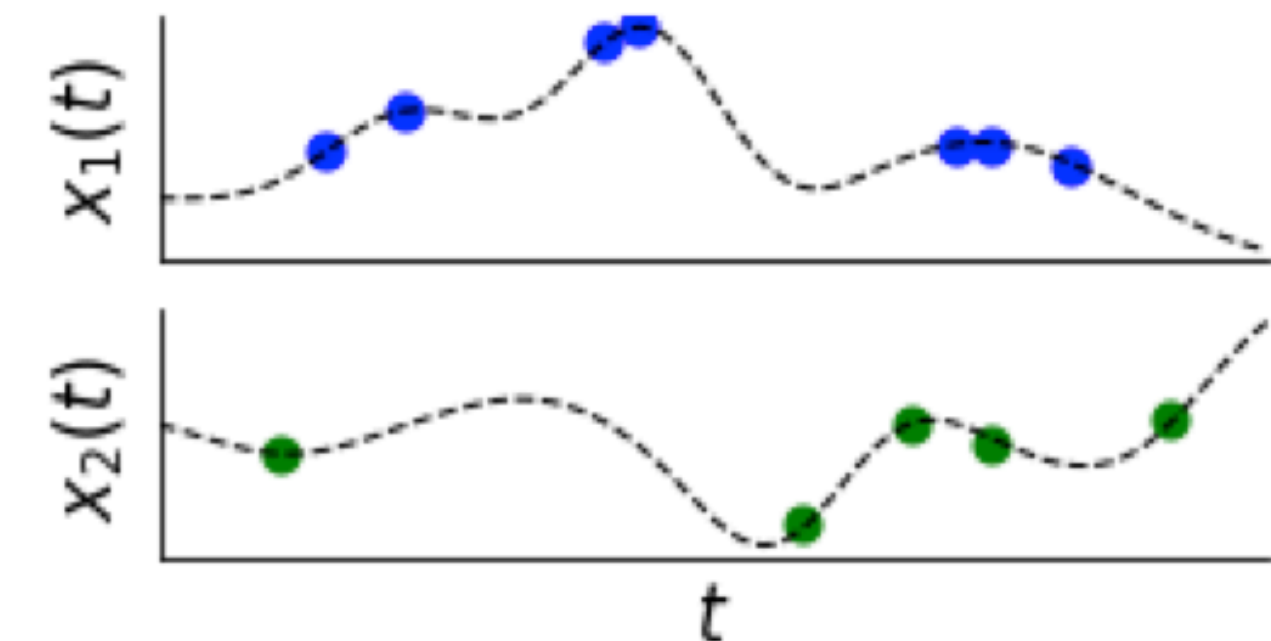
# Regularly and Irregularly Sampled Time Series



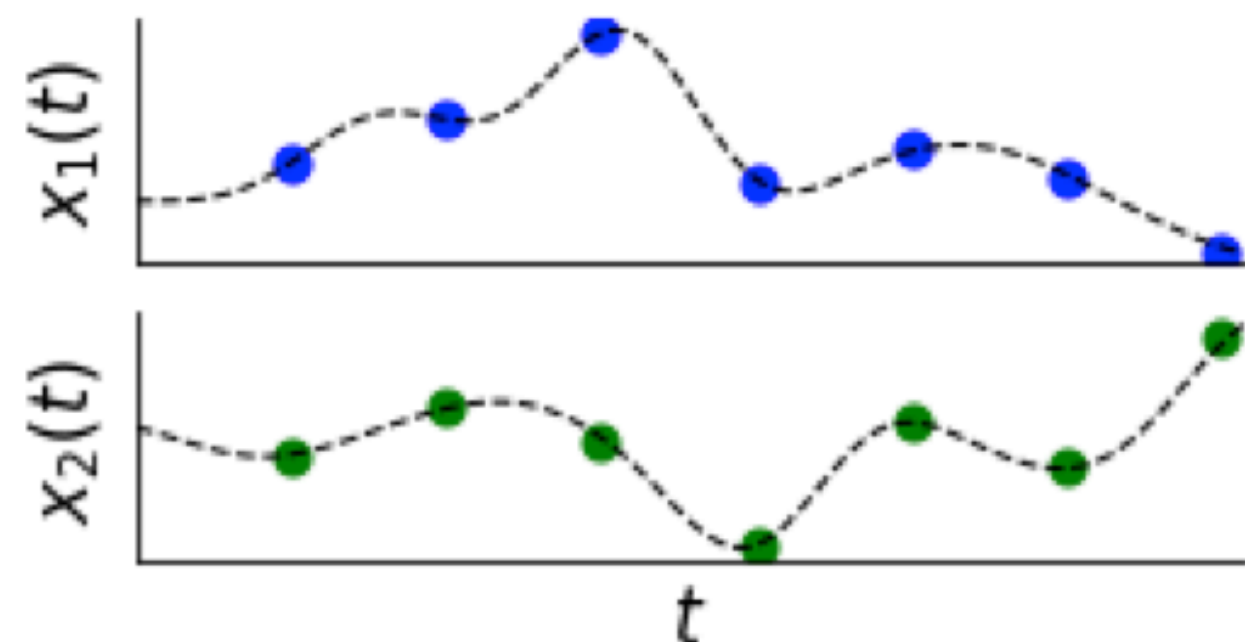
Univariate regularly sampled



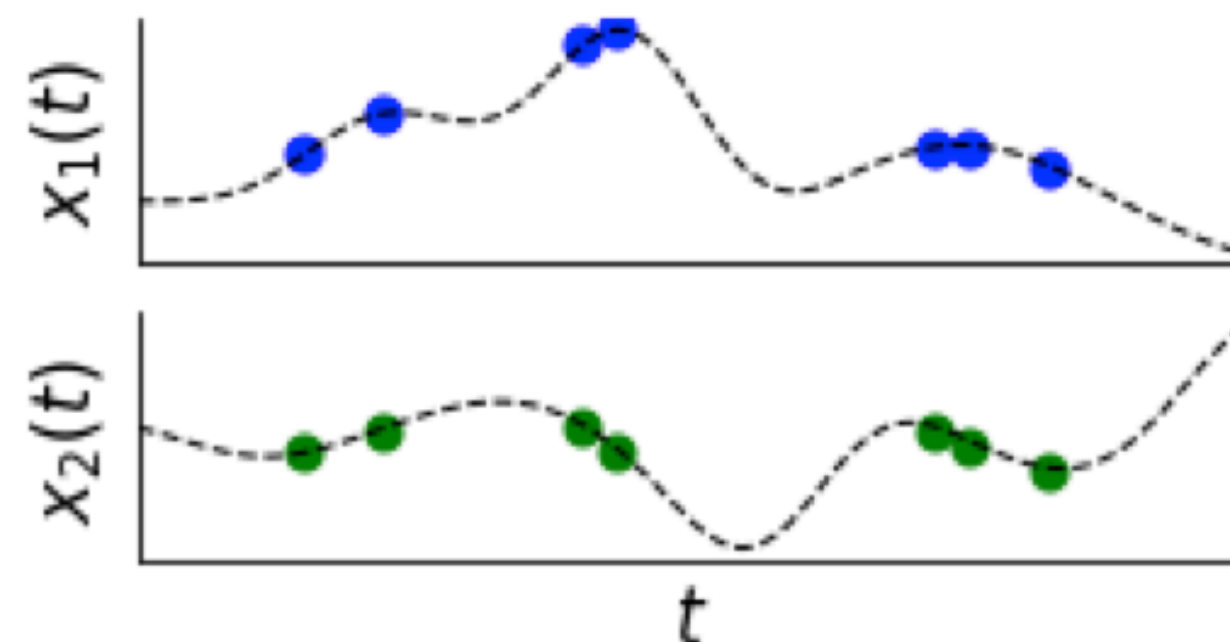
Univariate irregularly sampled



Multivariate irregularly sampled (unaligned)



Multivariate regularly sampled



Multivariate irregularly sampled (aligned)

Can we learn classification models on irregularly sampled time series without prior imputation?  
(imputation can be unreliable and sacrifice (which sacrifices interpretability))

# Why do we get Irregular Sampling?

---

# Why do we get Irregular Sampling?

---

- Availability of sensors
  - ▶ energy, mobility, multi-tenancy, human operator



# Why do we get Irregular Sampling?

---

- Availability of sensors
  - ▶ energy, mobility, multi-tenancy, human operator
- Samples are lost or dropped
  - ▶ network outage, bit corruption

# Why do we get Irregular Sampling?

---

- Availability of sensors
  - ▶ energy, mobility, multi-tenancy, human operator
- Samples are lost or dropped
  - ▶ network outage, bit corruption
- Sensors report asynchronous events
  - ▶ e.g. motion sensors, event-oriented imagers

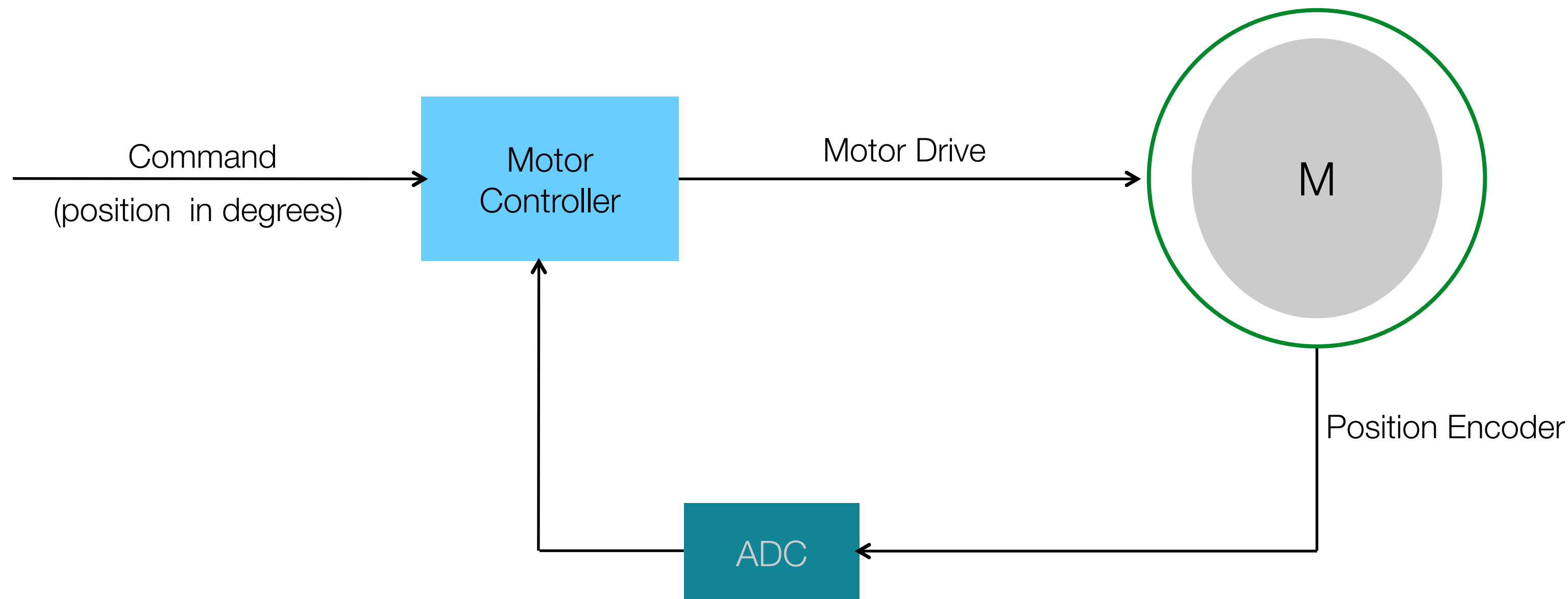
# Why do we get Irregular Sampling?

---

- Availability of sensors
  - ▶ energy, mobility, multi-tenancy, human operator
- Samples are lost or dropped
  - ▶ network outage, bit corruption
- Sensors report asynchronous events
  - ▶ e.g. motion sensors, event-oriented imagers
- Sensors sampling intervals are adapted
  - ▶ state of the process being sampled (e.g. patient's health)

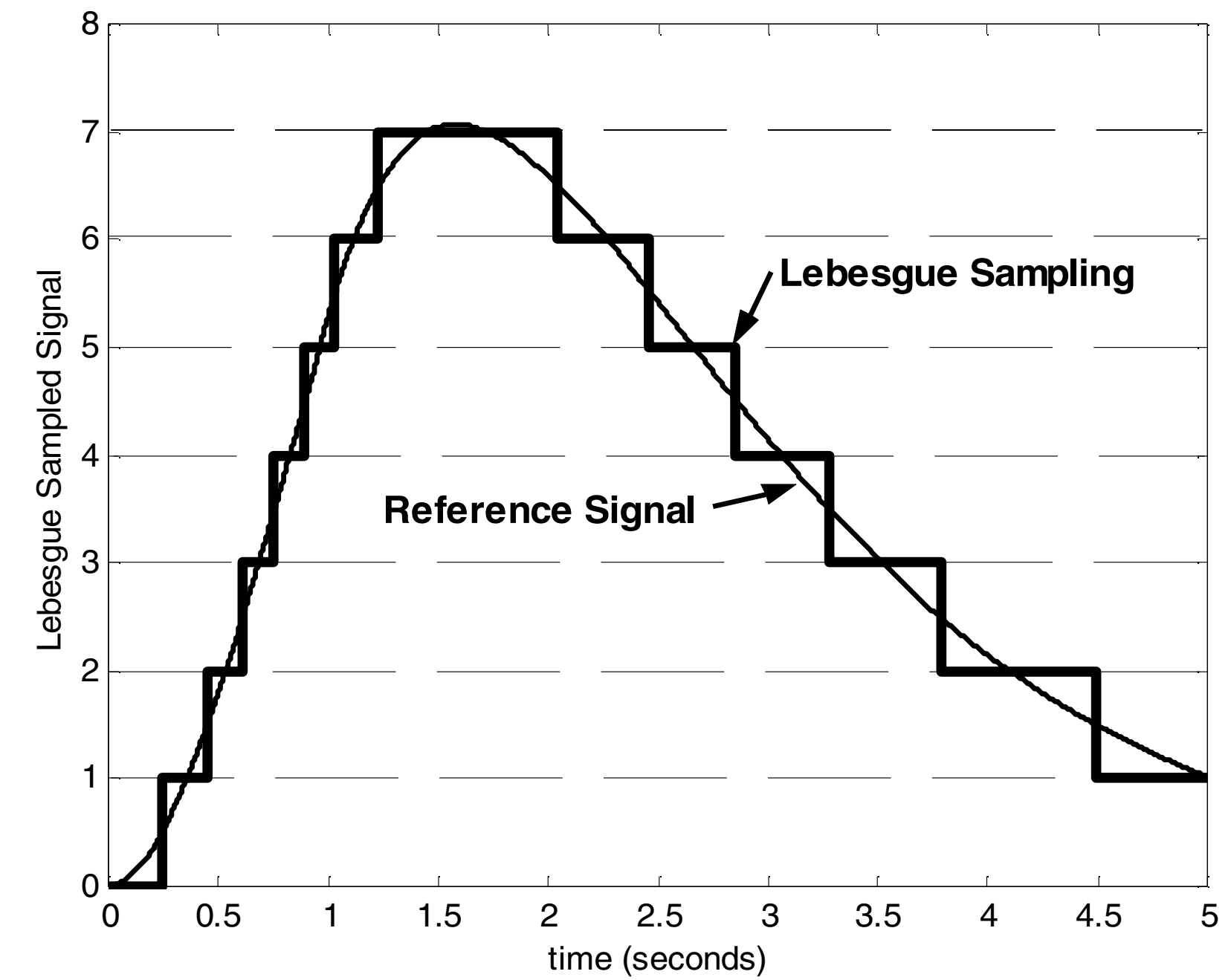
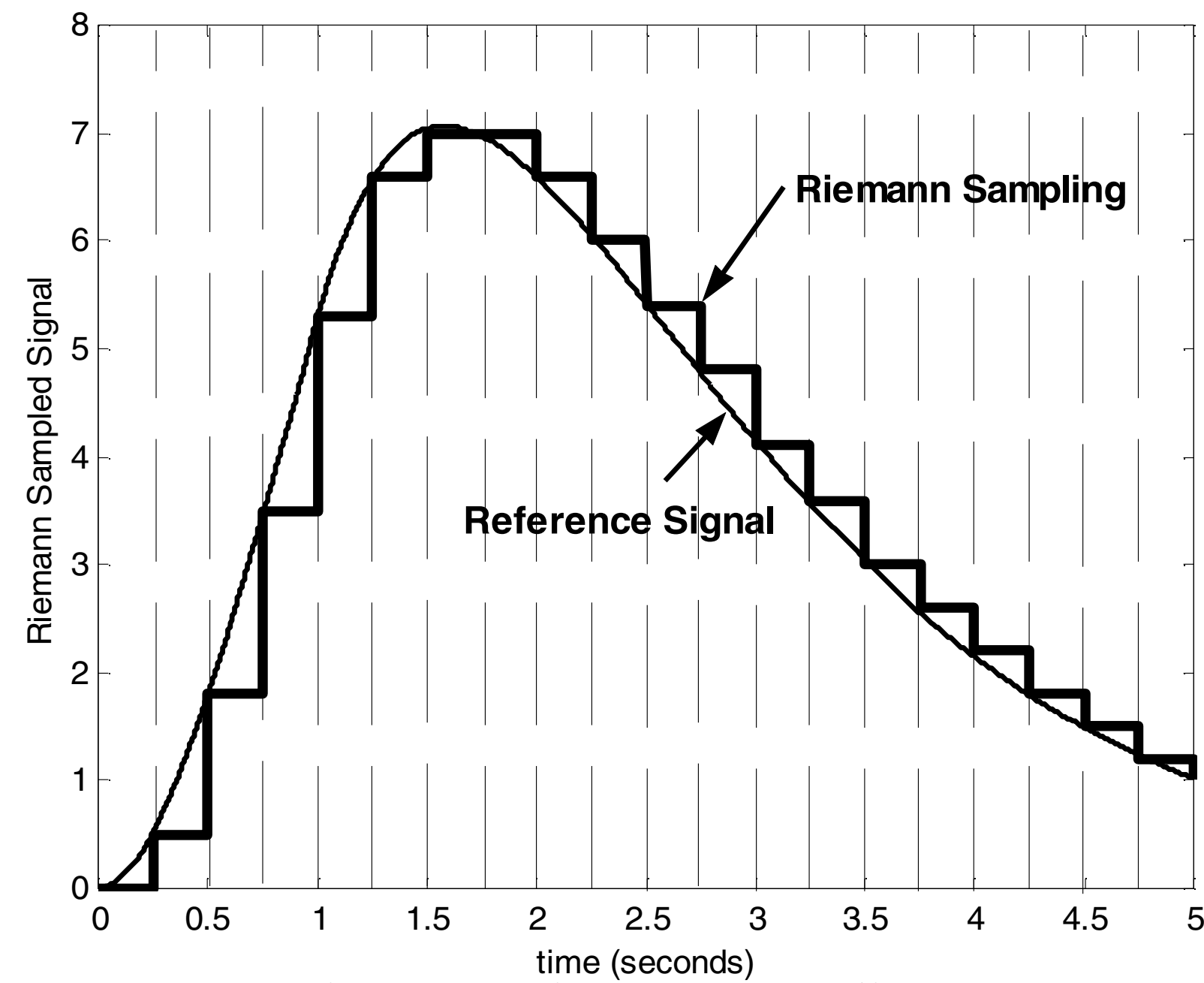
# Sampling Adaptively

---



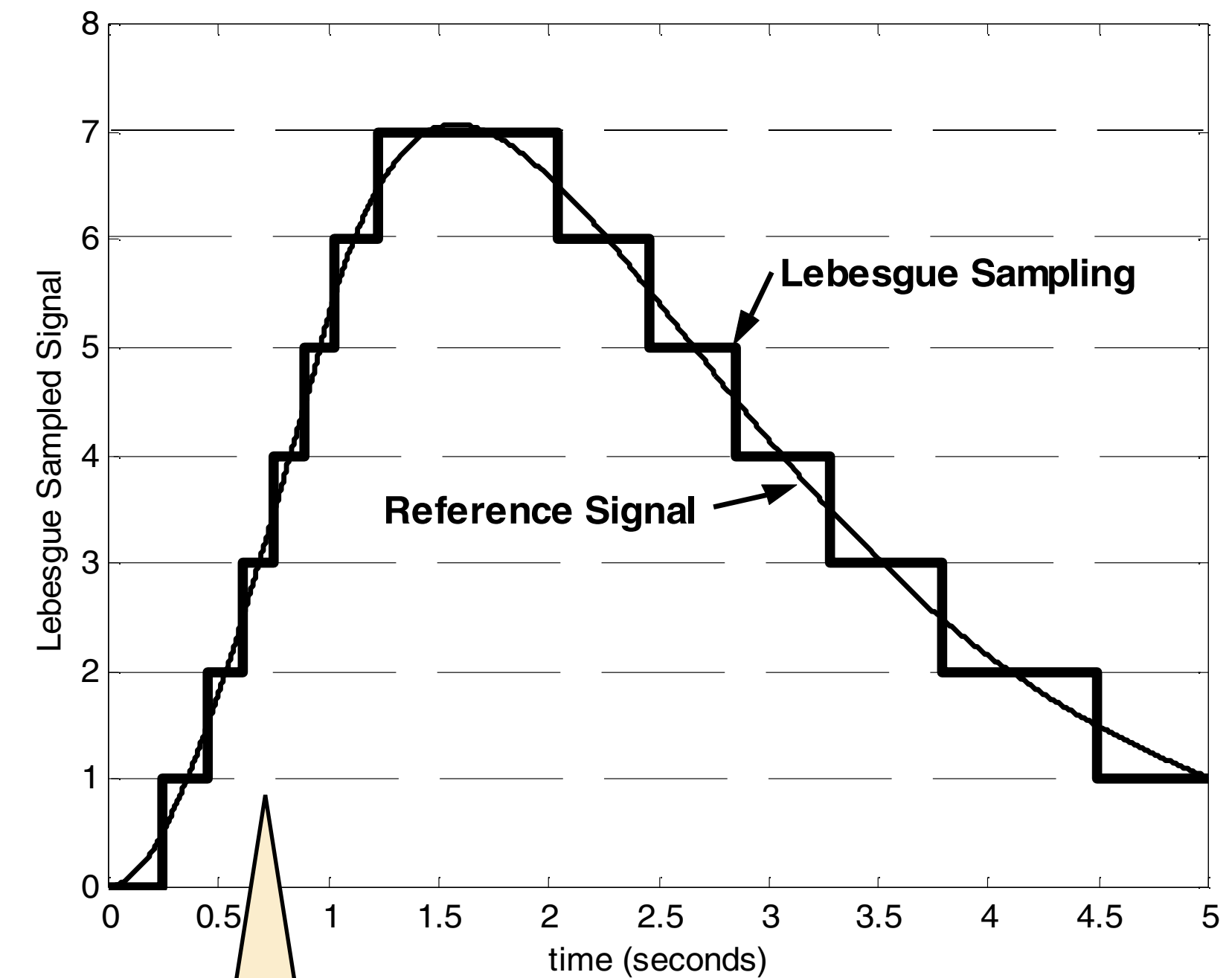
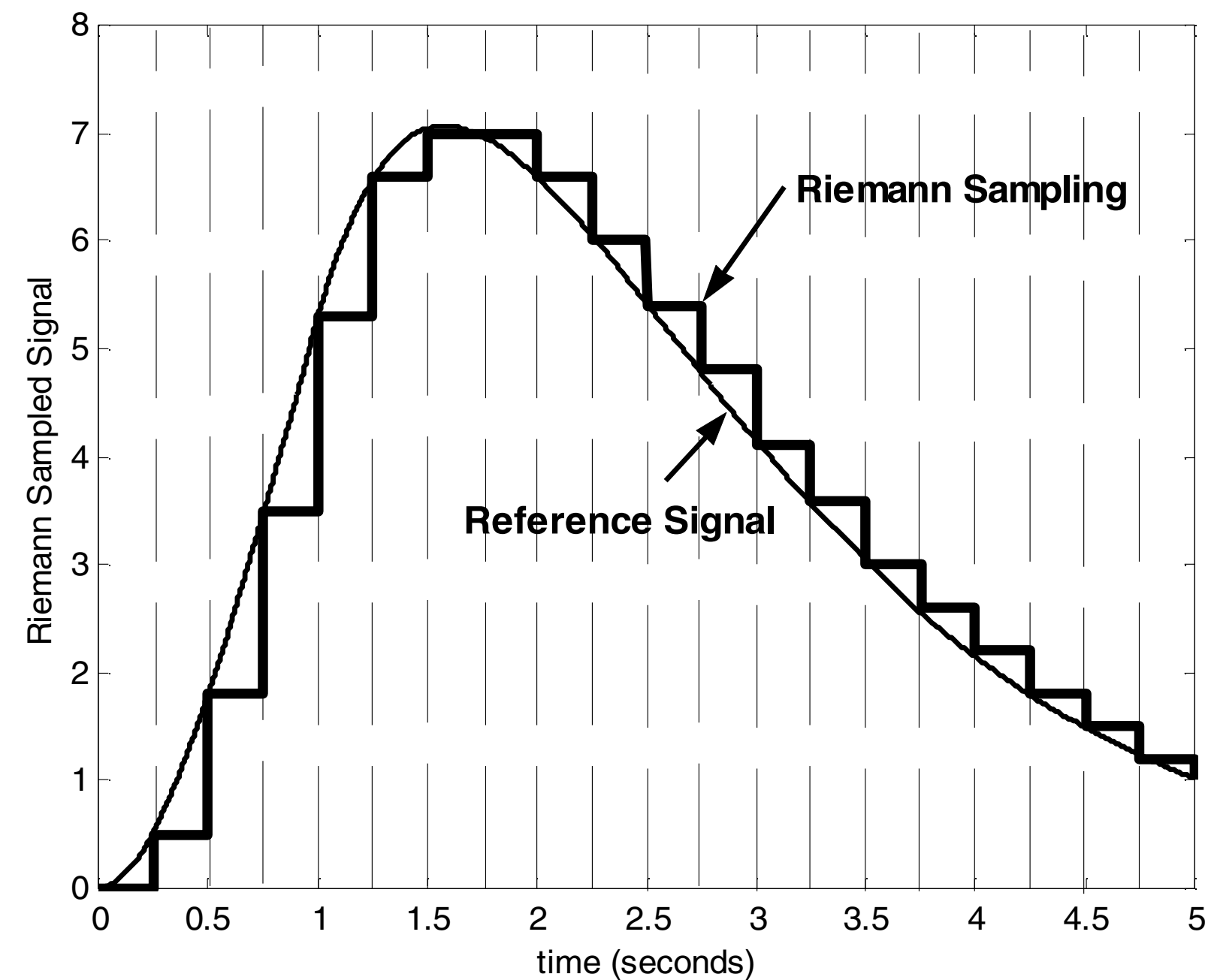
- Idea: Instead of sampling periodically as is traditionally done, could we sample only when the position of the motor changes sufficiently?
  - sample more frequently if position changes more rapidly
- This sampling policy is called *Lebesgue* Sampling

# Riemann Sampling vs. Lebesgue Sampling



From: Roy McCann, Anil Kumar Gunda, Suchit Reddy Damugatla, "Improved Operation of Networked Control Systems using Lebesgue Sampling"

# Riemann Sampling vs. Lebesgue Sampling



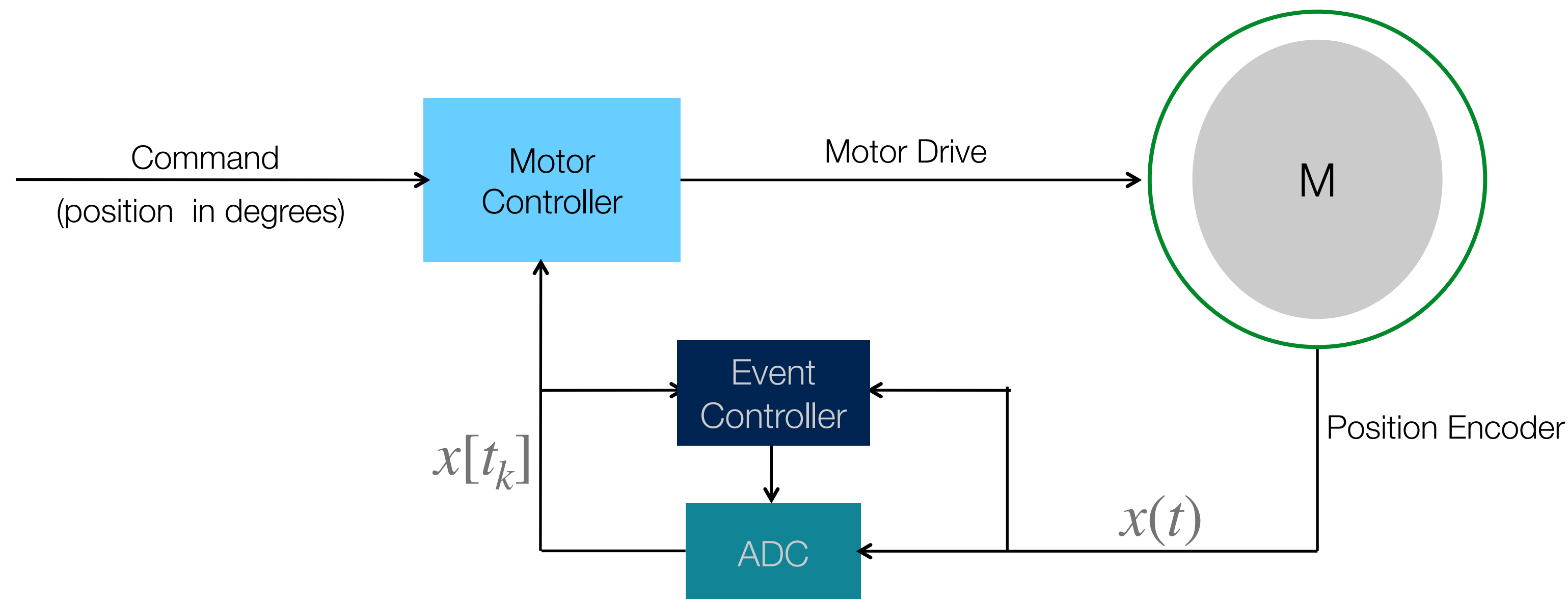
From: Roy McCann, An Improved  
Operation of Netw

How do we decide the Lebesgue interval?

Improved  
Sampling

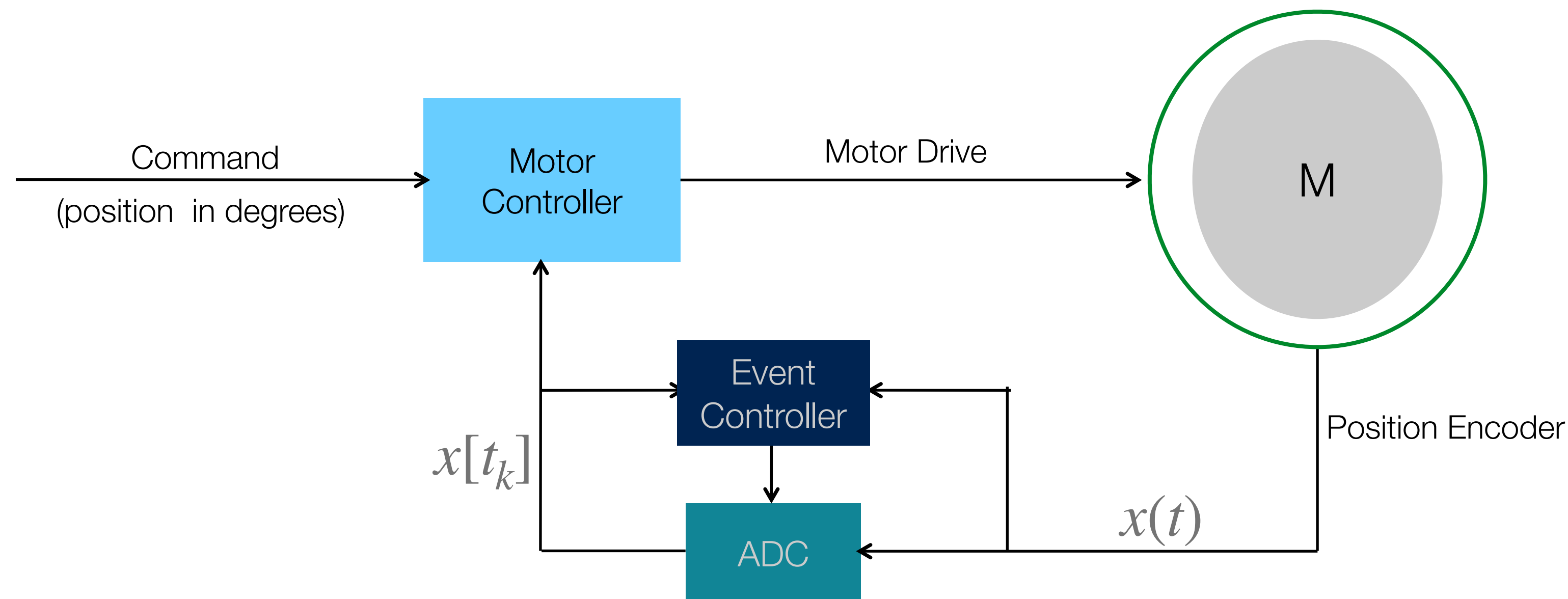
# Event-Triggered Control System

---



- Introduce an Event Controller that monitors the Position Encoder signal continuously and determines the optimal sampling time
- $x(t)$  is the continuous-time analog position encoder signal
- $x[t_k]$  is the last sample produced by the ADC

# Event-Triggered Control System



- **Theorem [Anta-Tabuada]:** Sampling is optimal and the system is guaranteed stability if the Event Controller produces a sampling trigger whenever:

$$||x(t) - x[t_k]|| > \sigma ||x(t)||$$

- ▶  $\sigma$  is a design parameter that trades off average sampling rate and performance
- ▶ Intuition: Lebesgue interval must be the “error” relative to size of signal



# Self-Triggered Control System

---

- Event Controller needs to check the inequality at ALL times. Can one do better?
- Motor controller is designed with a model of the motor in mind
- Given the current state of the motor and the command input, control theory provides an estimate of the next state
- One can extend this estimation to predict when the state of the motor will violate the sampling inequality
- Thus, one can predict when the sample should be taken and schedule the ADC

$$\textit{CurrentMeasurement} + \textit{SystemDynamics} \Rightarrow \textit{NextSampleTime}$$

- No continuous time check needed
- However, this does not work well when the model has uncertainties or when perturbations to the system cannot be bounded

# Self-Triggered Control System

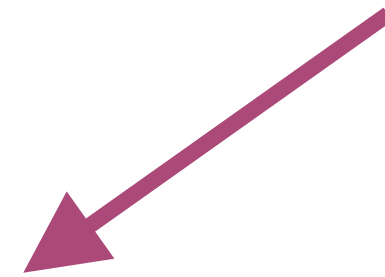
---

$$t_{k+1} = \alpha(\lambda, x[t_k])$$

# Self-Triggered Control System

---

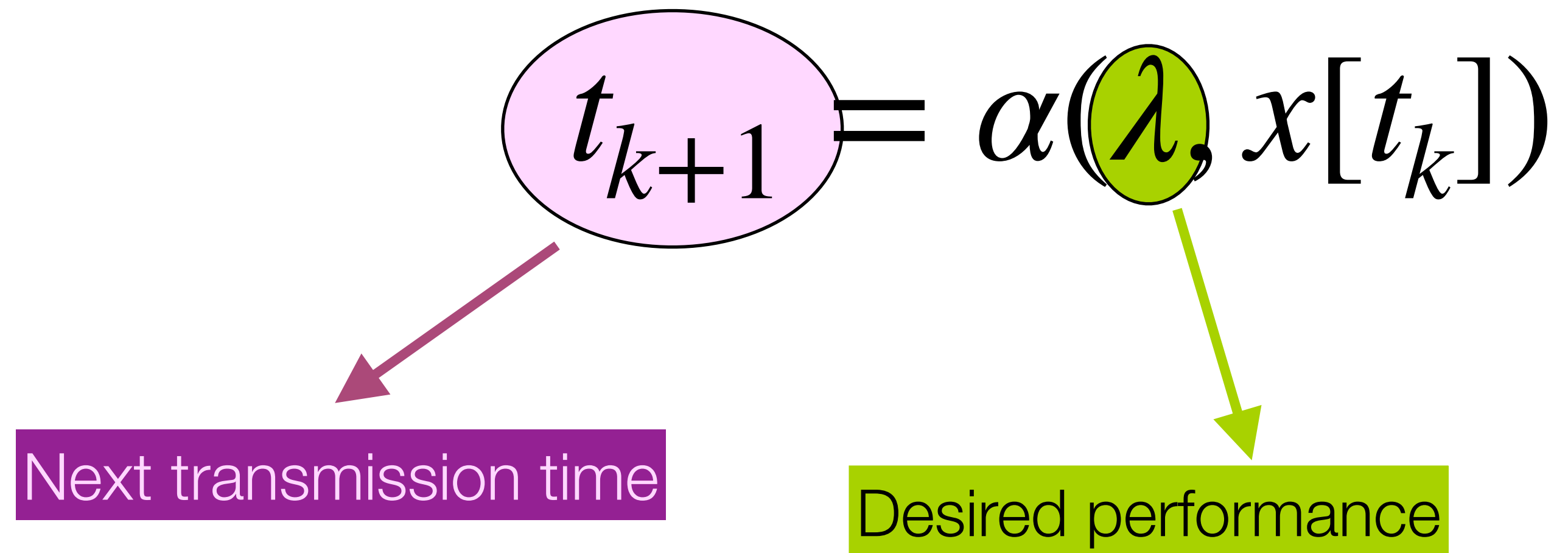
$$t_{k+1} = \alpha(\lambda, x[t_k])$$



Next transmission time

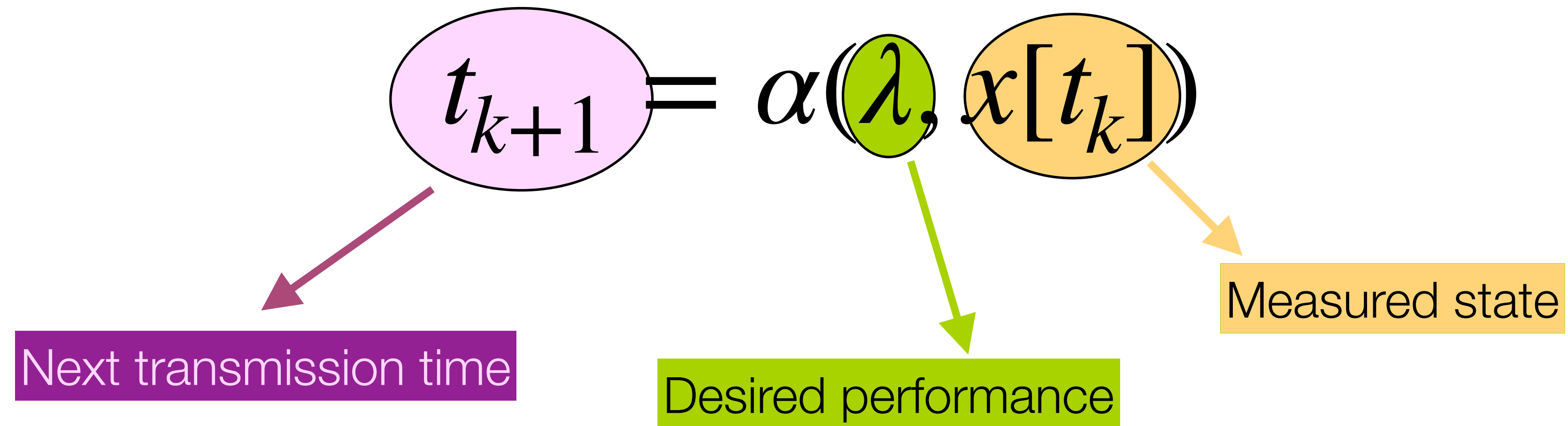
# Self-Triggered Control System

---



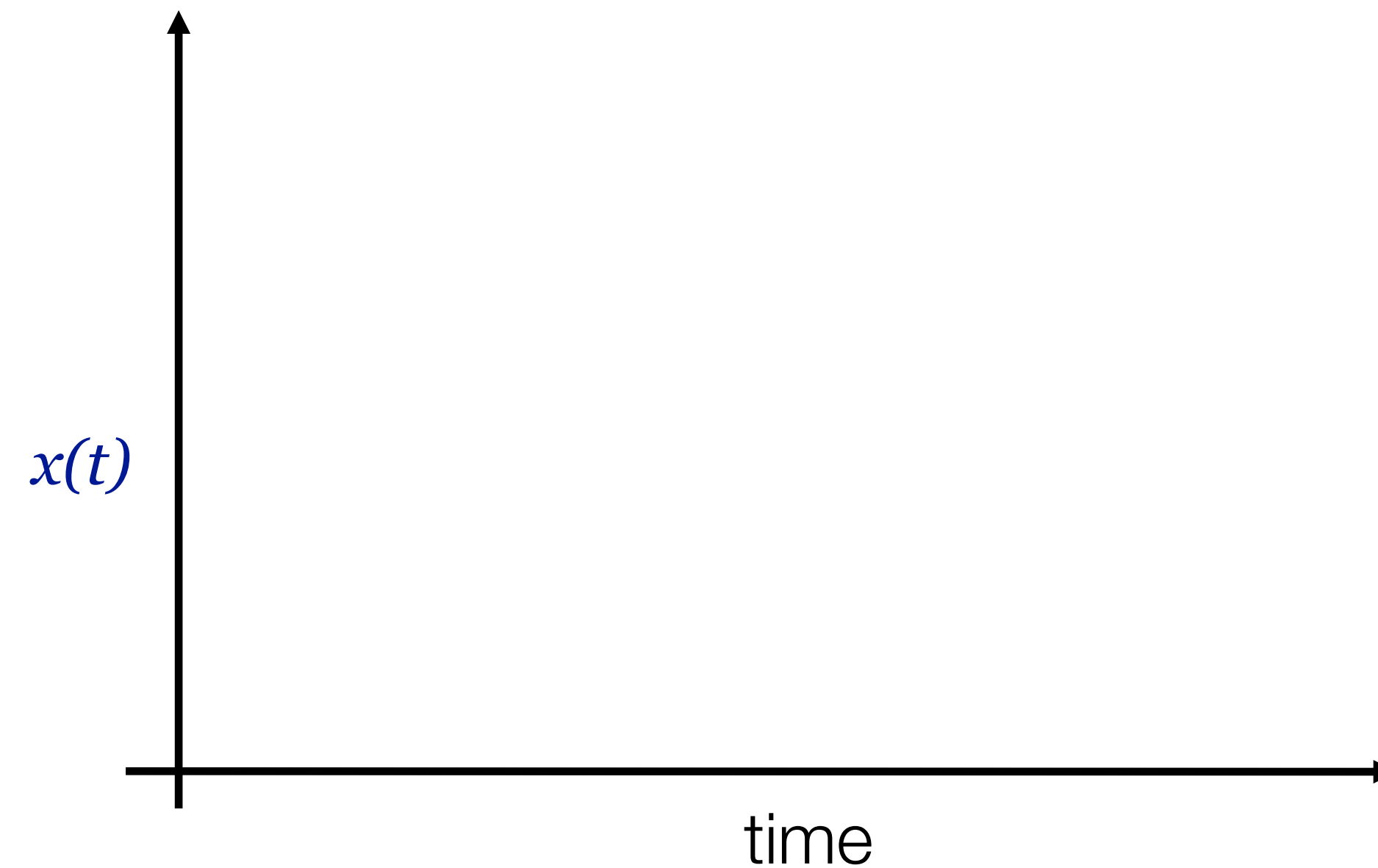
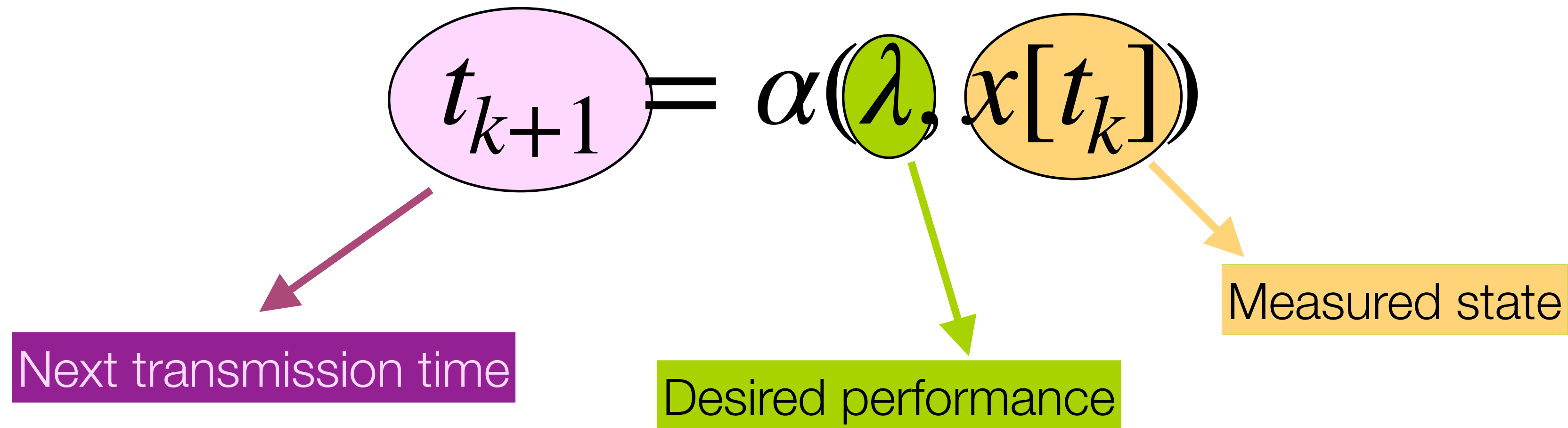
# Self-Triggered Control System

---



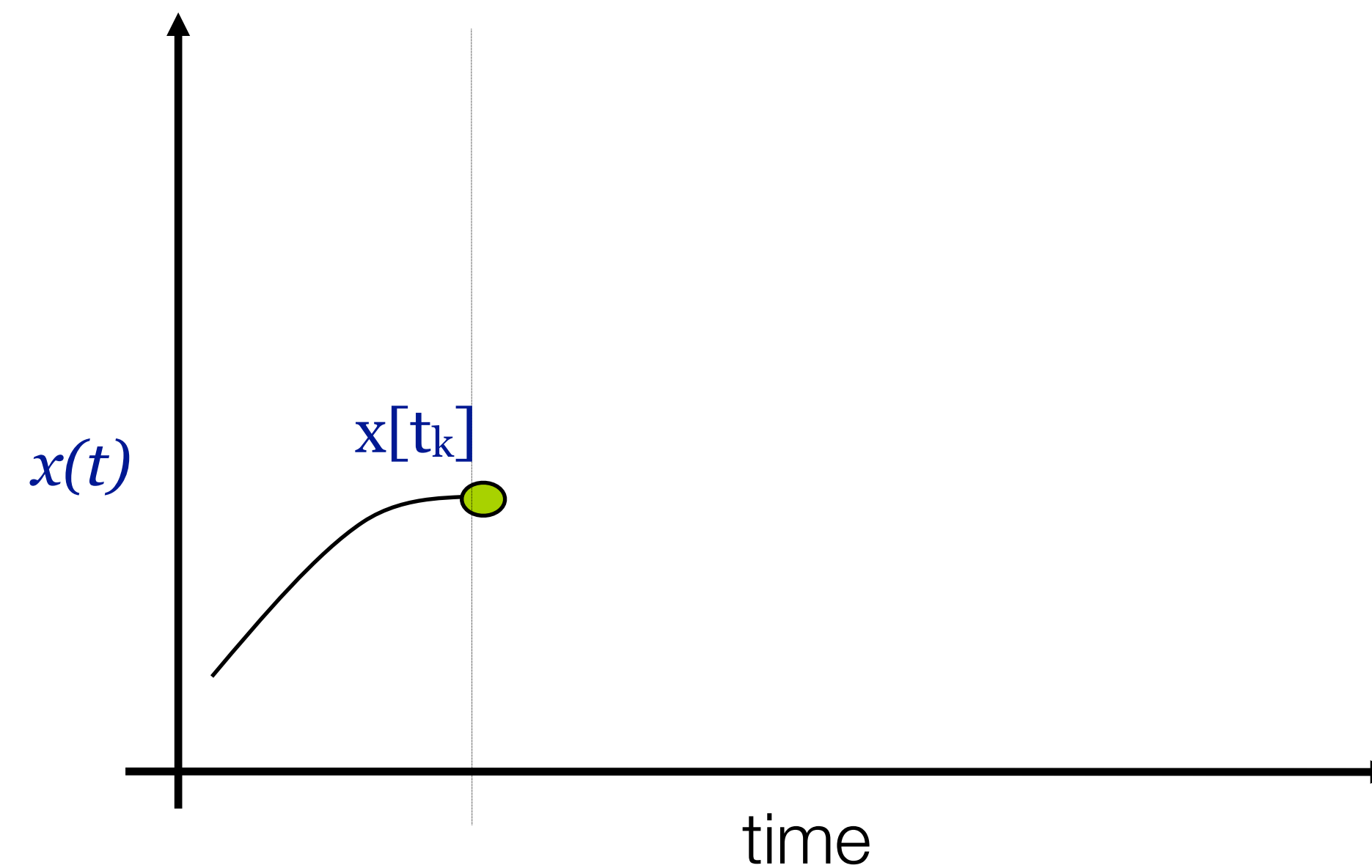
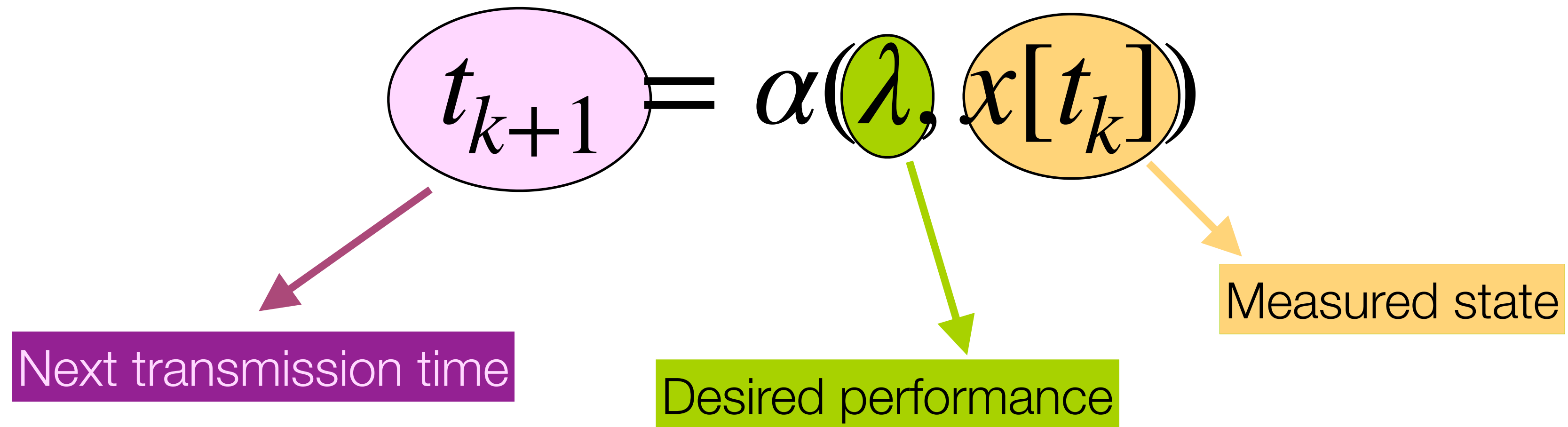
# Self-Triggered Control System

---



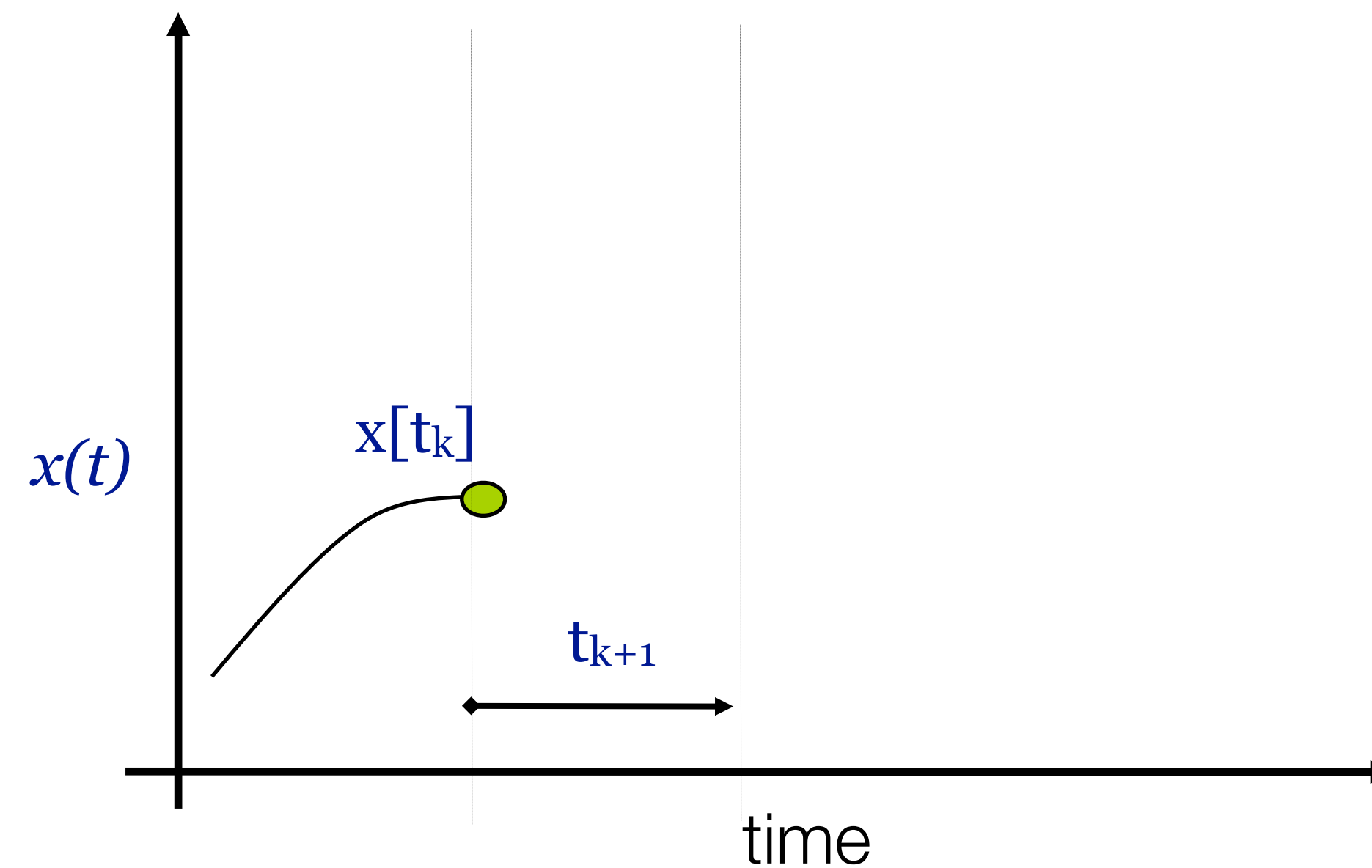
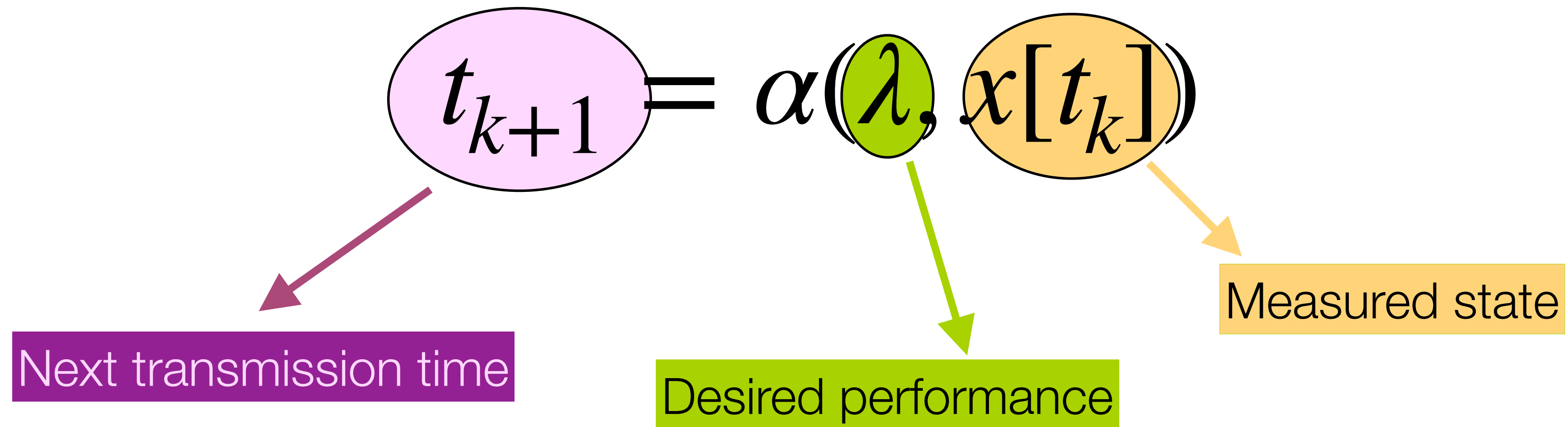
# Self-Triggered Control System

---



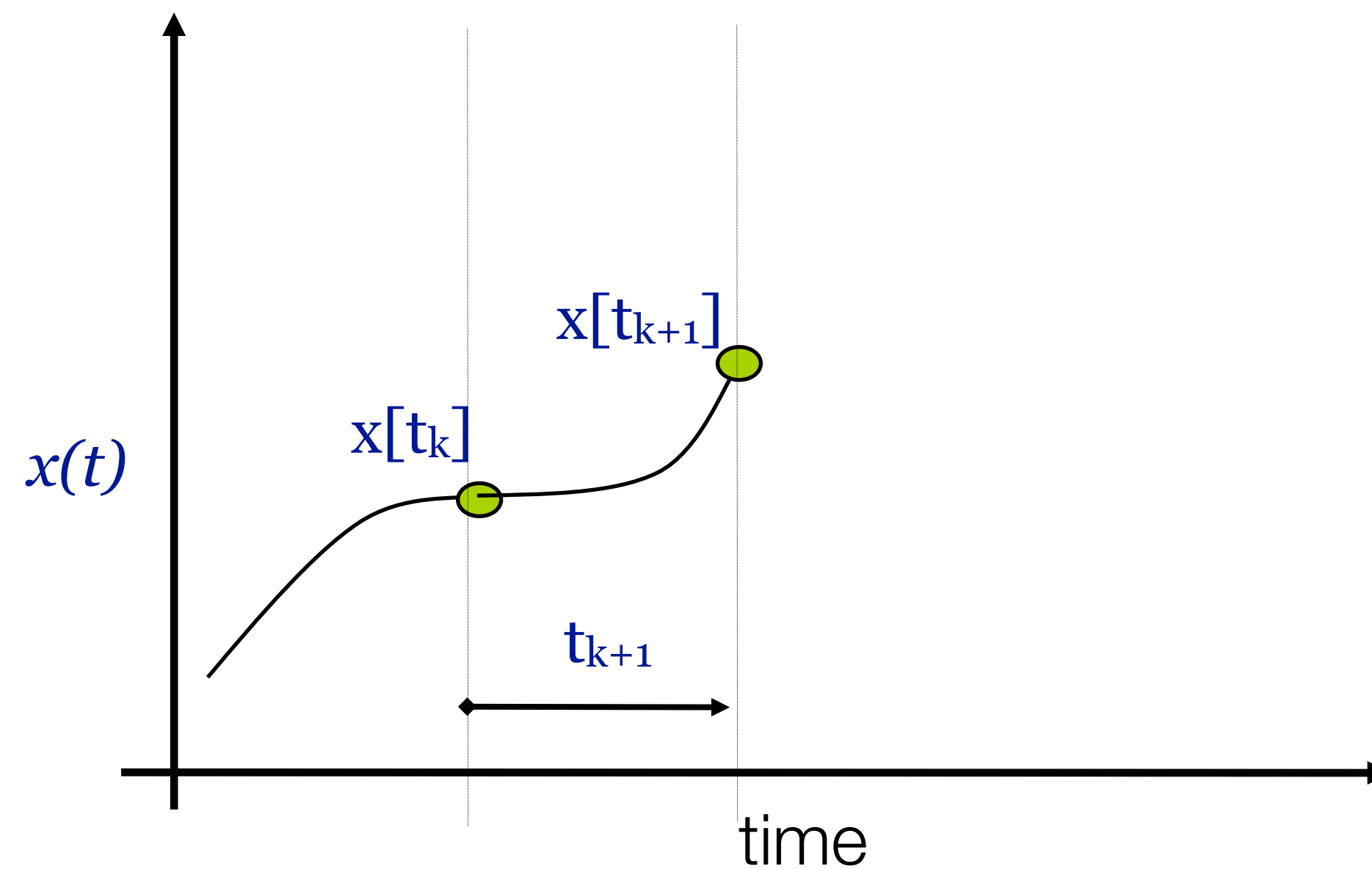
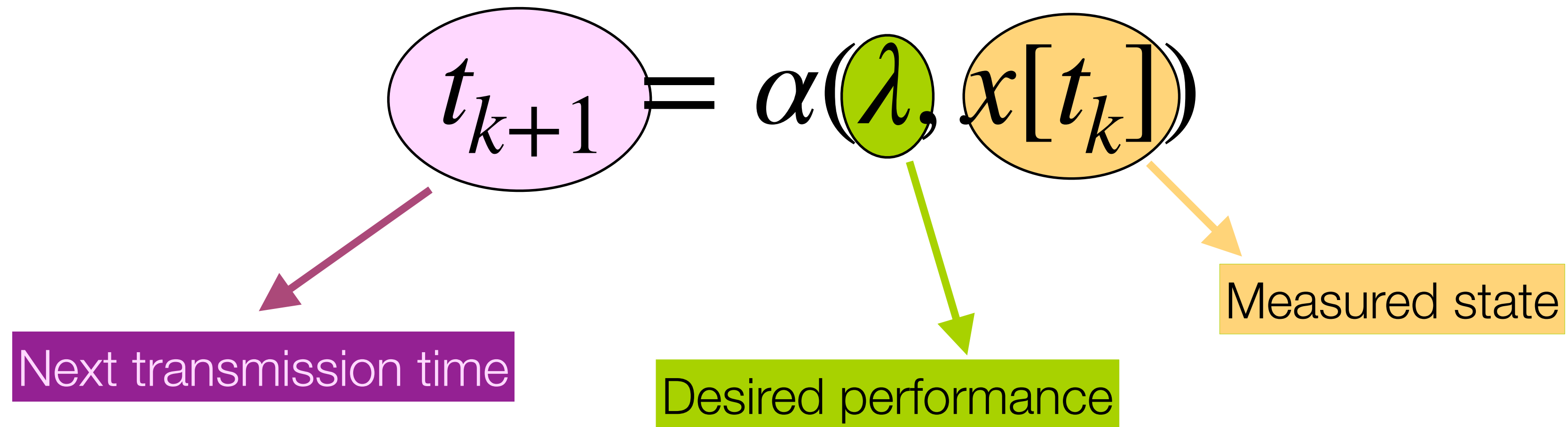
# Self-Triggered Control System

---

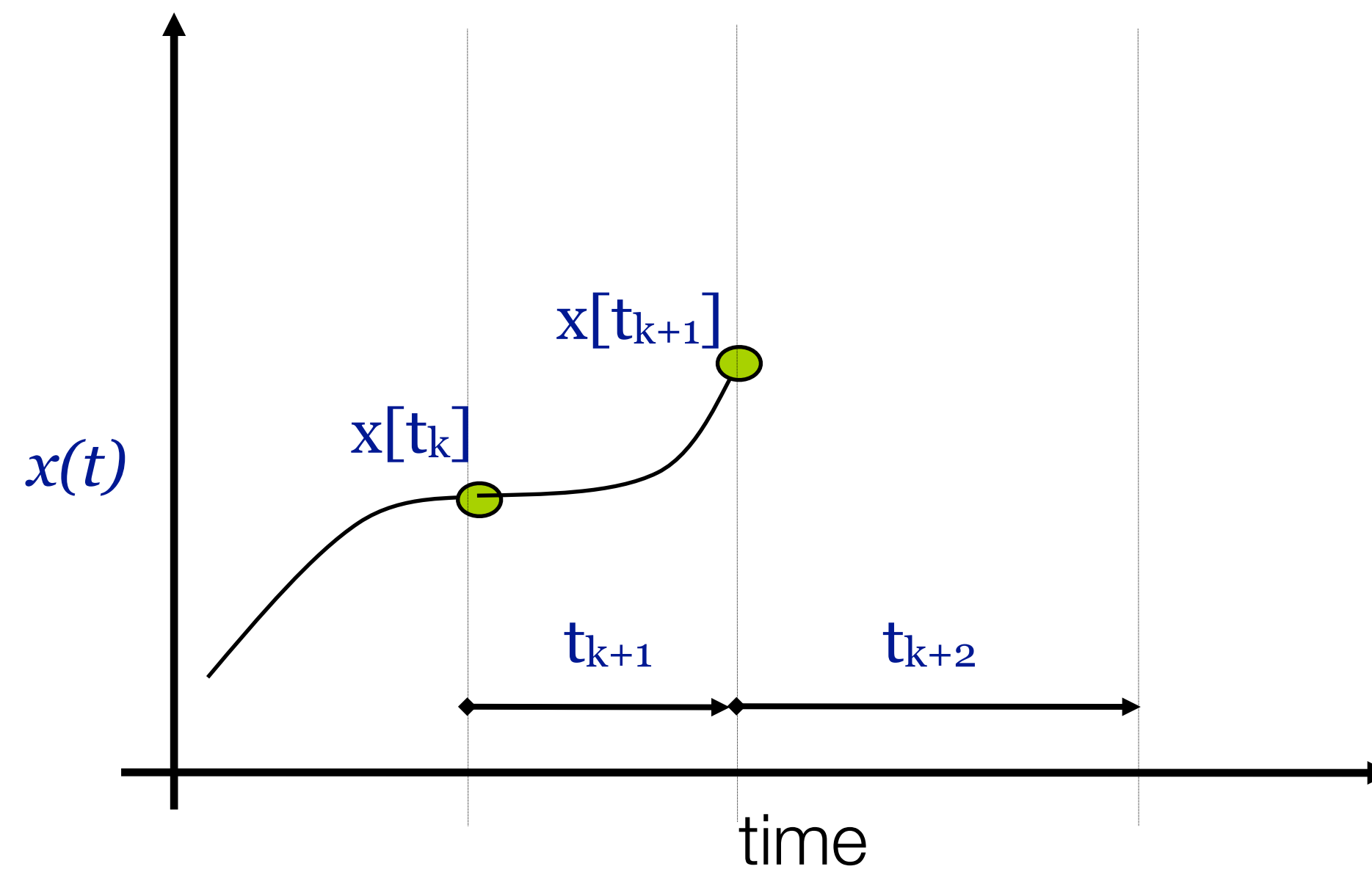
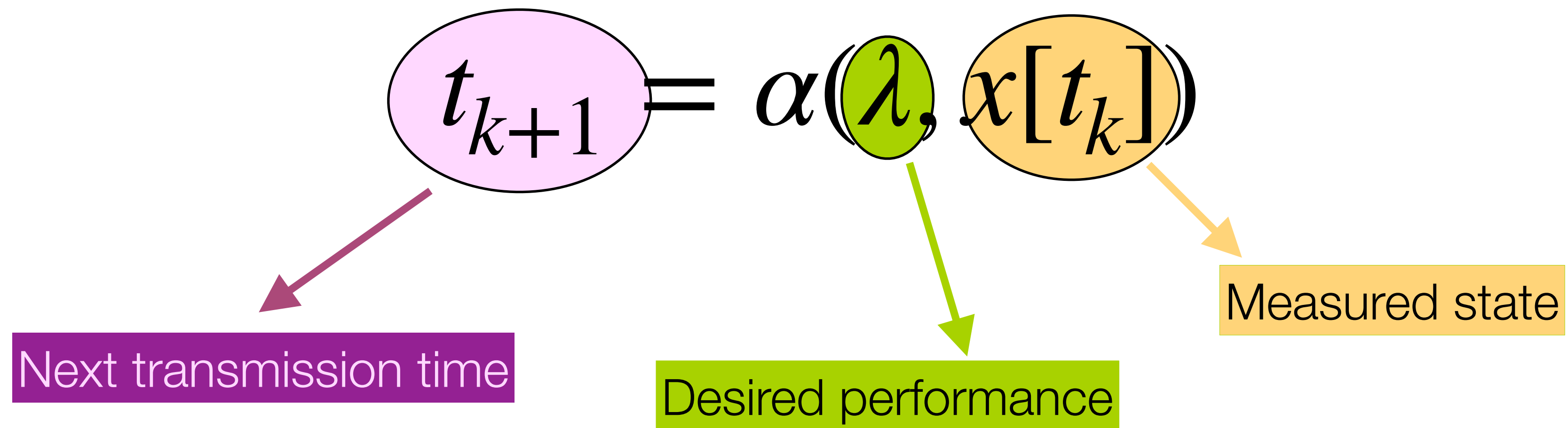




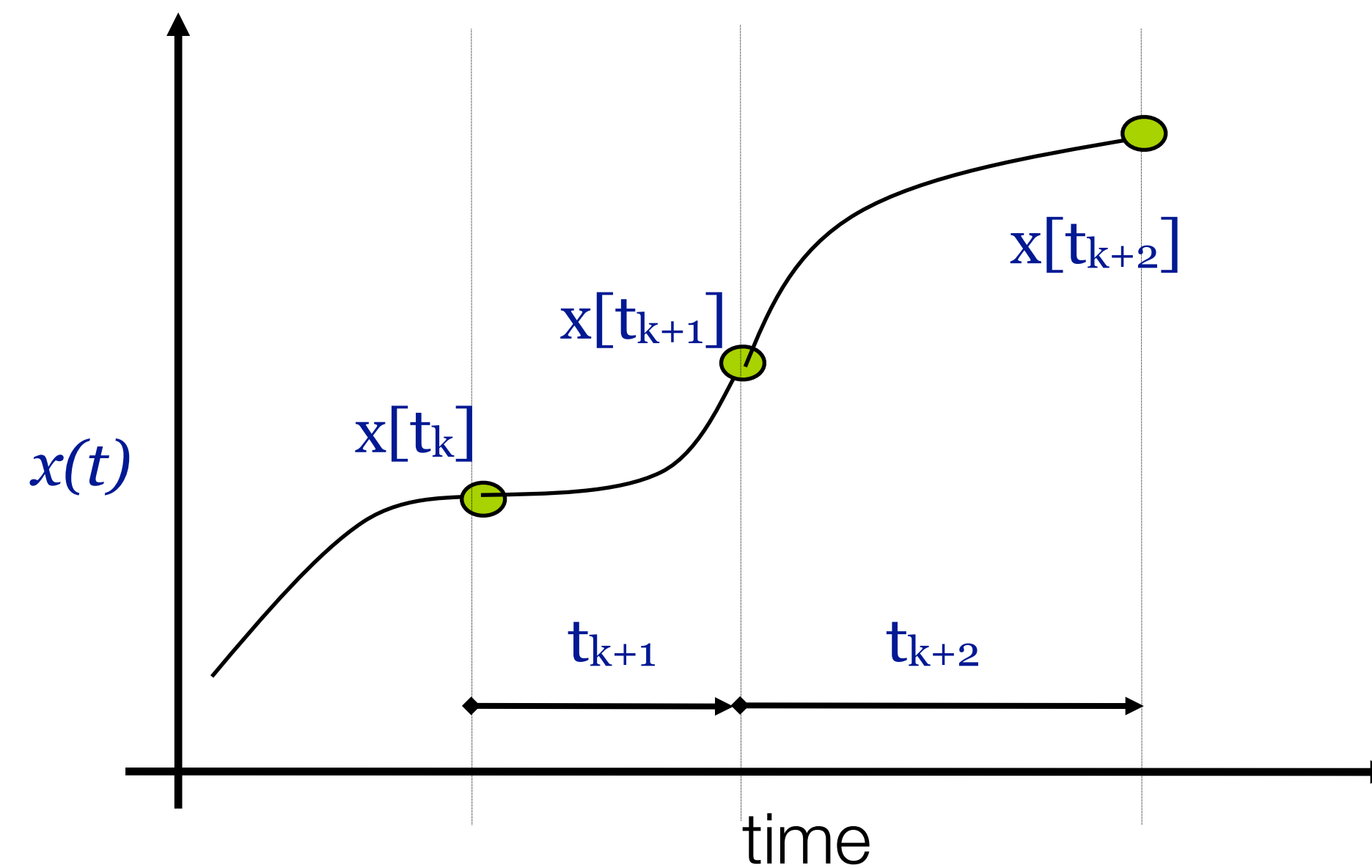
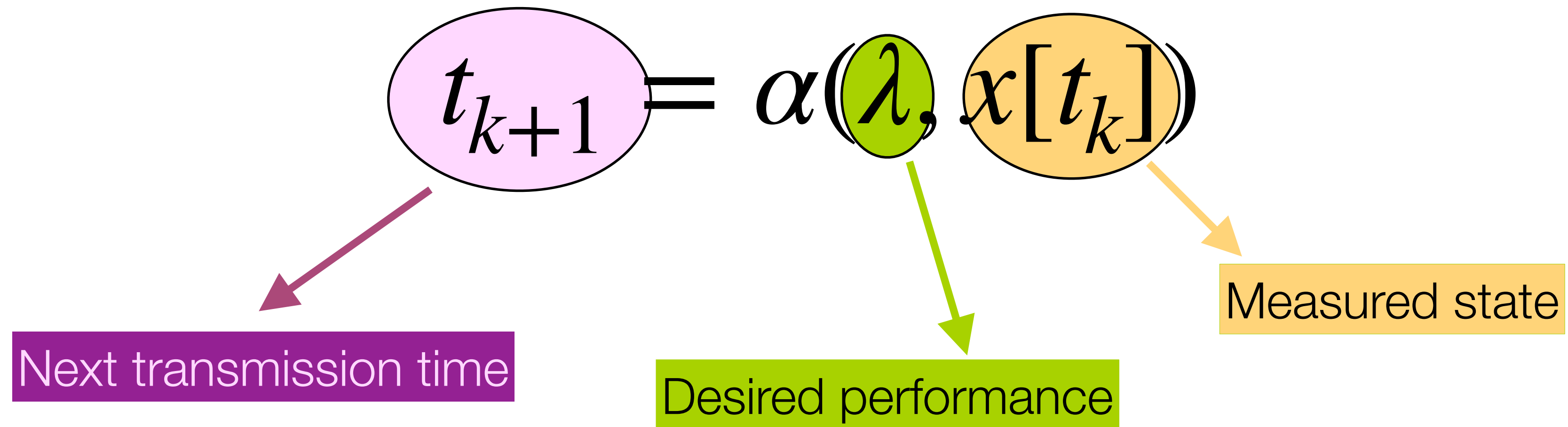
# Self-Triggered Control System



# Self-Triggered Control System



# Self-Triggered Control System

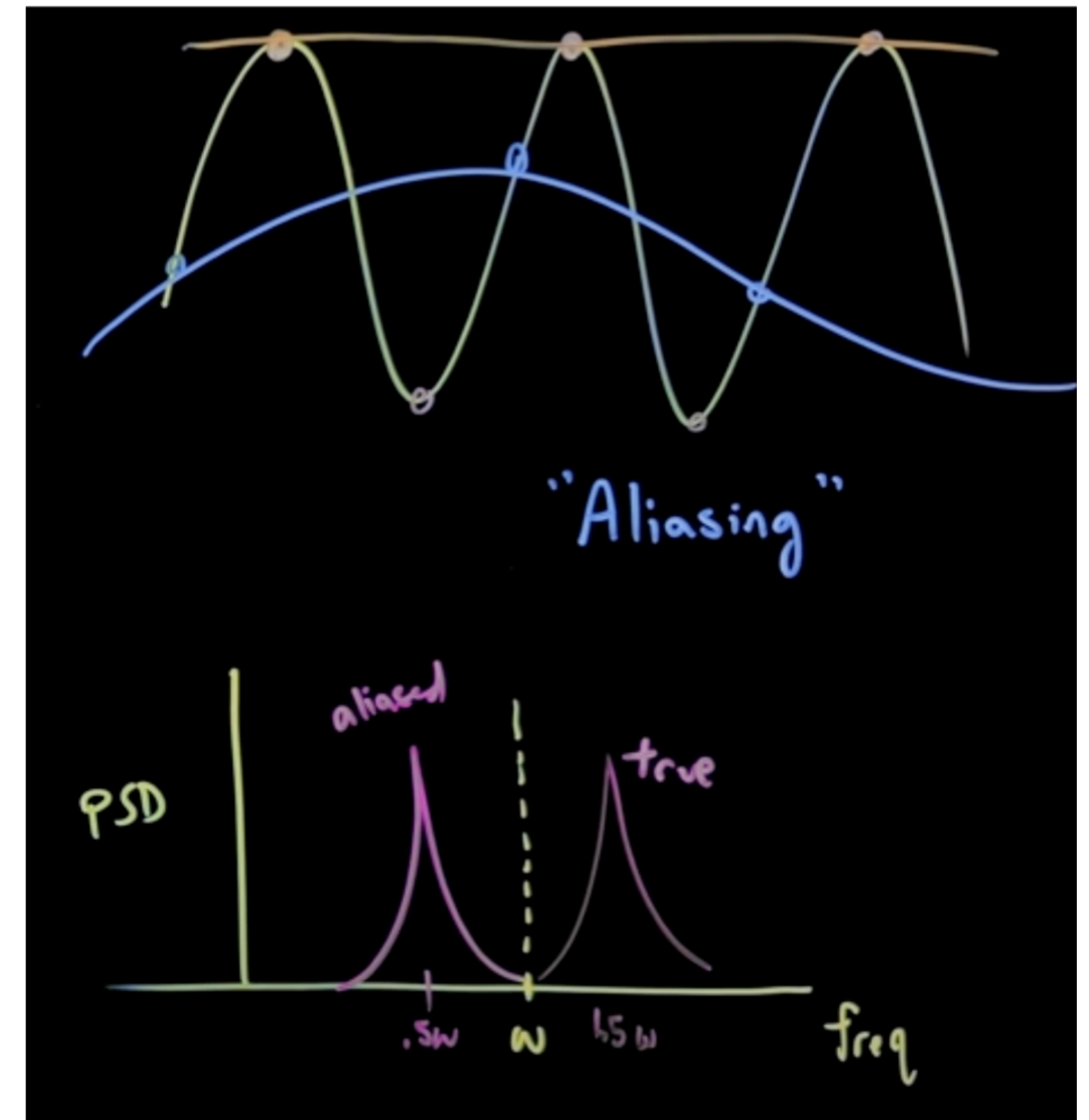


# Why do we get Irregular Sampling?

---

- Availability of sensors
  - ▶ energy, mobility, multi-tenancy, human operator
- Samples are lost or dropped
  - ▶ network outage, bit corruption
- Sensors report asynchronous events
  - ▶ e.g. motion sensors, event-oriented imagers
- Sensors sampling intervals are adapted
  - ▶ state of the process being sampled (e.g. patient's health)
- Compressive sampling of sensors

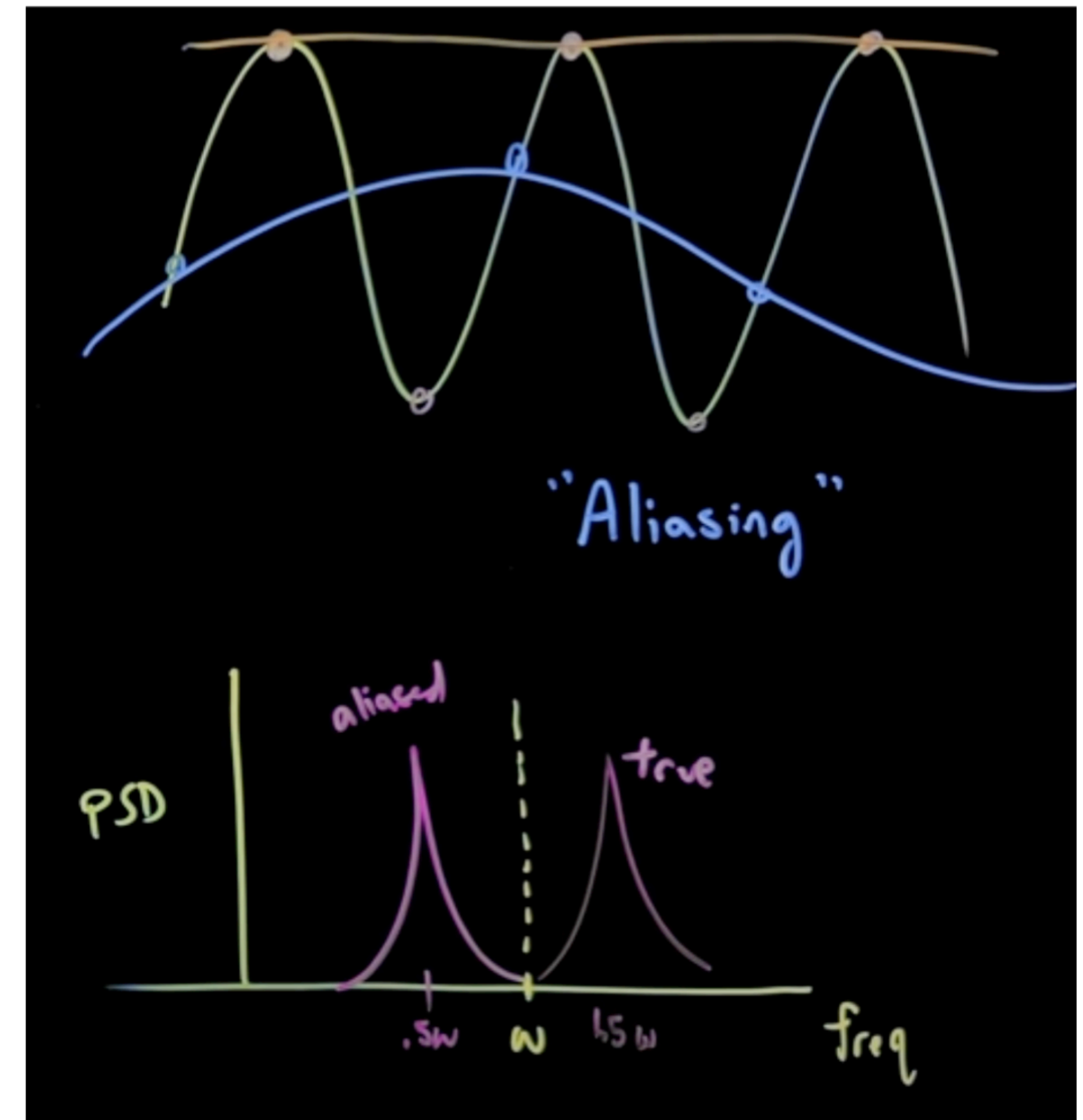
# Sampling a Signal



Steve Brunton, "Shannon Nyquist Sampling Theorem"  
<https://www.youtube.com/watch?v=FcXZ28BX-xE&t=449s>

# Sampling a Signal

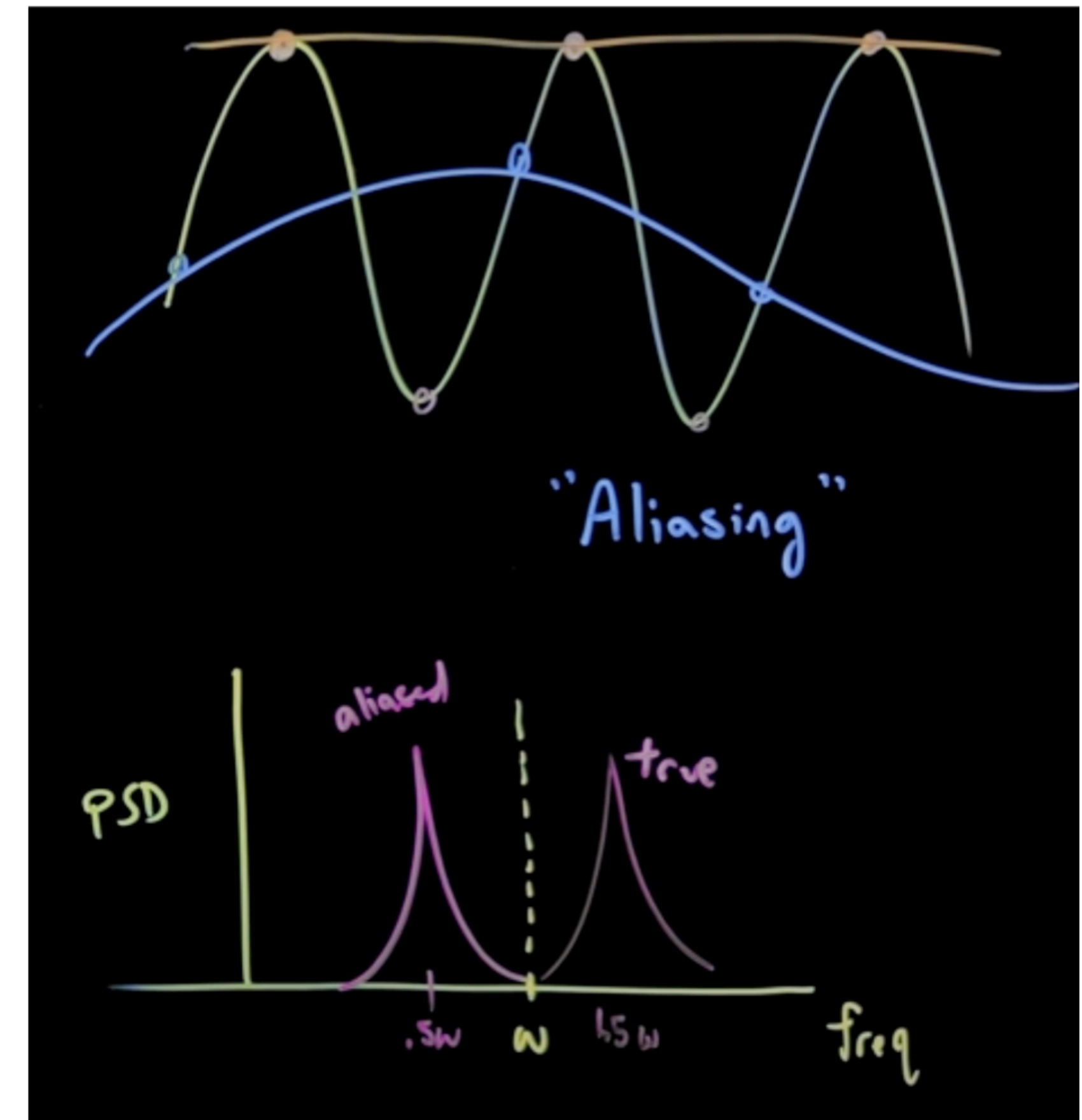
- Shannon-Nyquist Sampling Theorem:
  - ▶ To resolve all frequencies in a function, it must be sampled at twice the highest frequency present
  - ▶ A function containing no frequency  $> \omega$  Hz is completely determined by sampling at  $2\omega$  Hz (Nyquist Rate)  $\Delta t = \frac{1}{2\omega}$
  - ▶ Aliasing if we sample at a rate lower than  $2\omega$



Steve Brunton, "Shannon Nyquist Sampling Theorem"  
<https://www.youtube.com/watch?v=FcXZ28BX-xE&t=449s>

# Sampling a Signal

- Shannon-Nyquist Sampling Theorem:
  - ▶ To resolve all frequencies in a function, it must be sampled at twice the highest frequency present
  - ▶ A function containing no frequency  $> \omega$  Hz is completely determined by sampling at  $2\omega$  Hz (Nyquist Rate)  $\Delta t = \frac{1}{2\omega}$
  - ▶ Aliasing if we sample at a rate lower than  $2\omega$
- Beating Shannon-Nyquist Sampling Theorem
  - ▶ Advances in applied mathematics, statistics, and optimization have changed how we think about sampling
  - ▶ Technically, the Shannon-Nyquist Sampling Theorem is necessary only for signals that are broadband, i.e. densely packed with energy in all the frequencies from 0 to  $\omega$
  - ▶ But if the signal is sparse in frequency domain, one can beat the  $2\omega$  sampling rate requirement

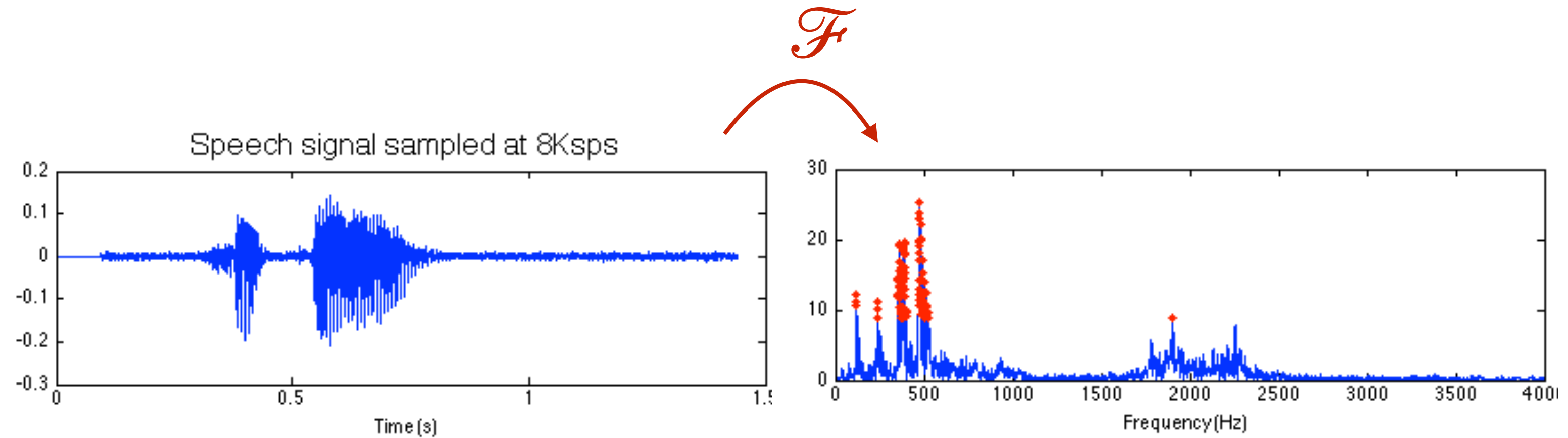
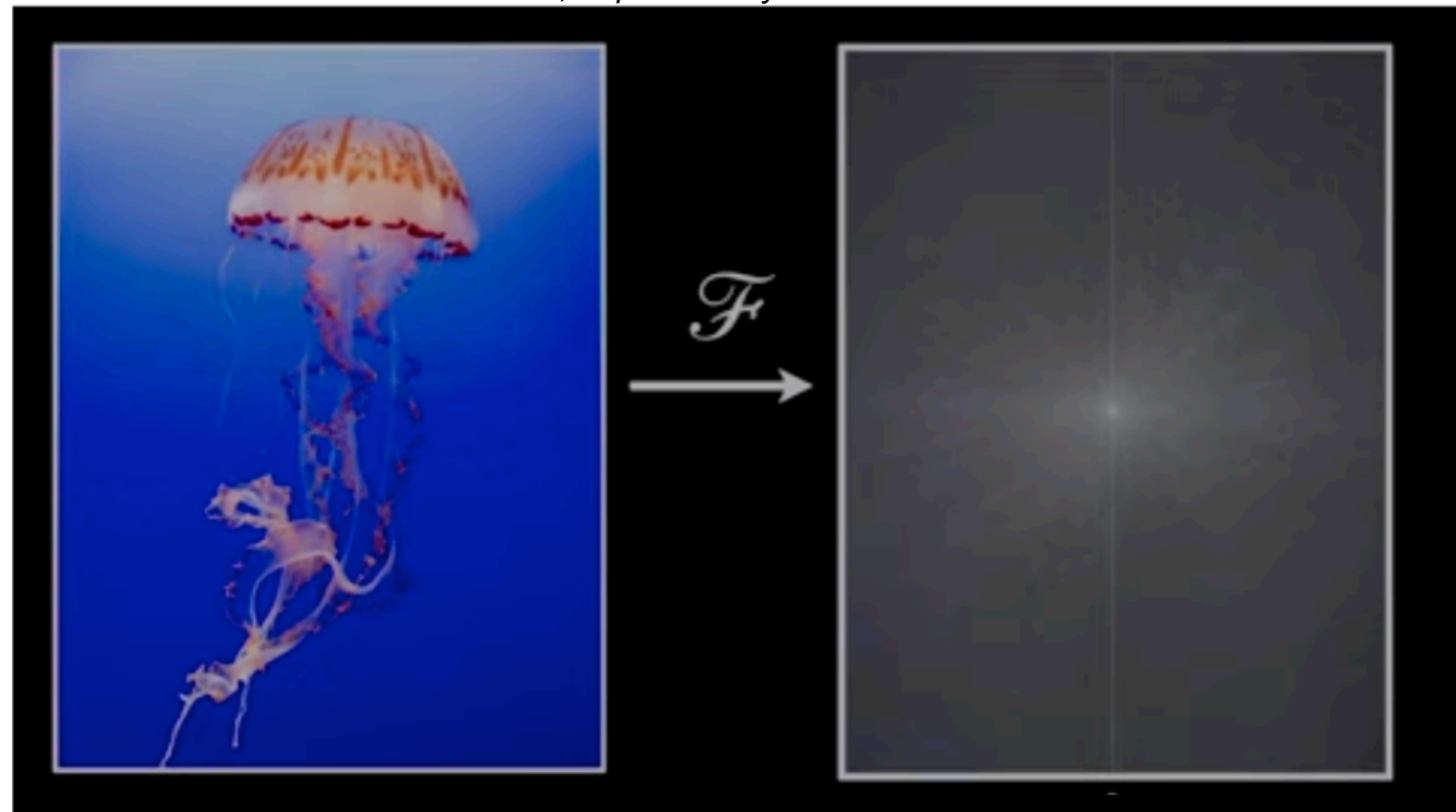


Steve Brunton, "Shannon Nyquist Sampling Theorem"  
<https://www.youtube.com/watch?v=FcXZ28BX-xE&t=449s>



# Motivating Compressive Sampling

Steve Brunton, <https://www.youtube.com/watch?v=hmBTACBGWJs>

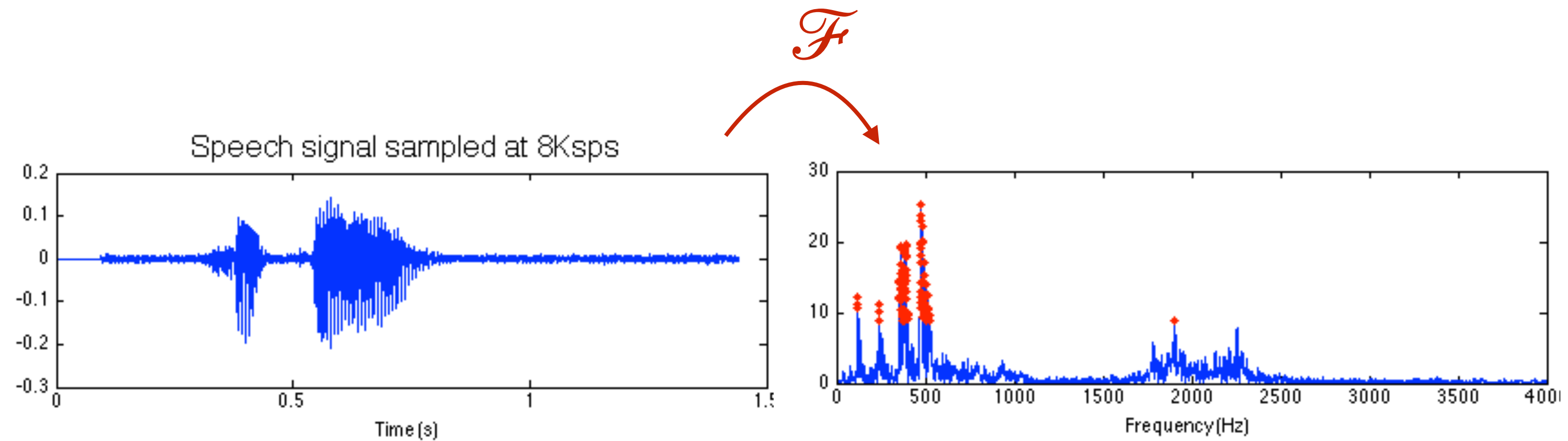
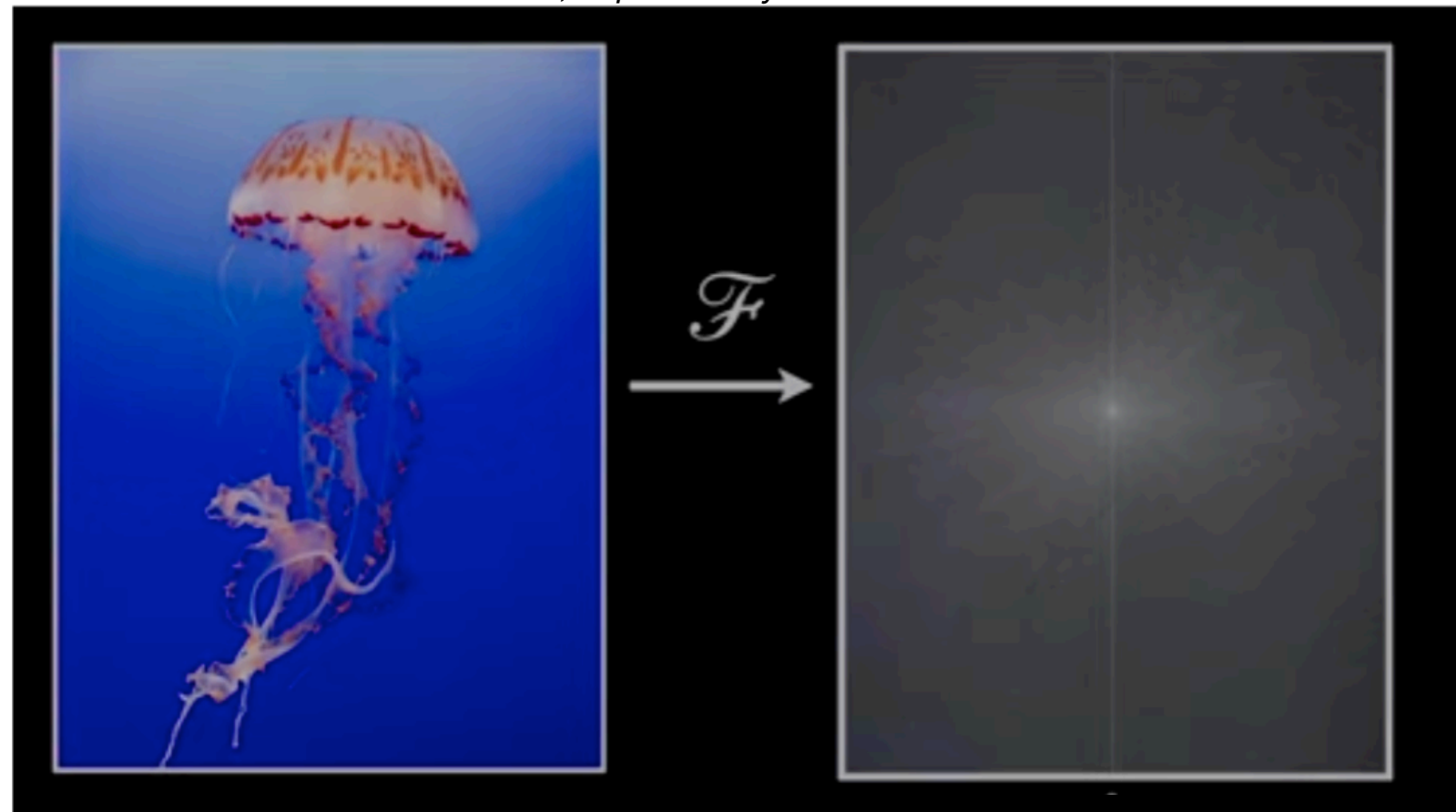


Real world signals are sparse in Fourier domain  
(or some other such universal domain, such as Wavelet)



# Motivating Compressive Sampling

Steve Brunton, <https://www.youtube.com/watch?v=hmBTACBGWJs>



Real world signals are sparse in Fourier domain  
(or some other such universal domain, such as Wavelet)

$$\mathbf{x} = \Psi \mathbf{s}$$

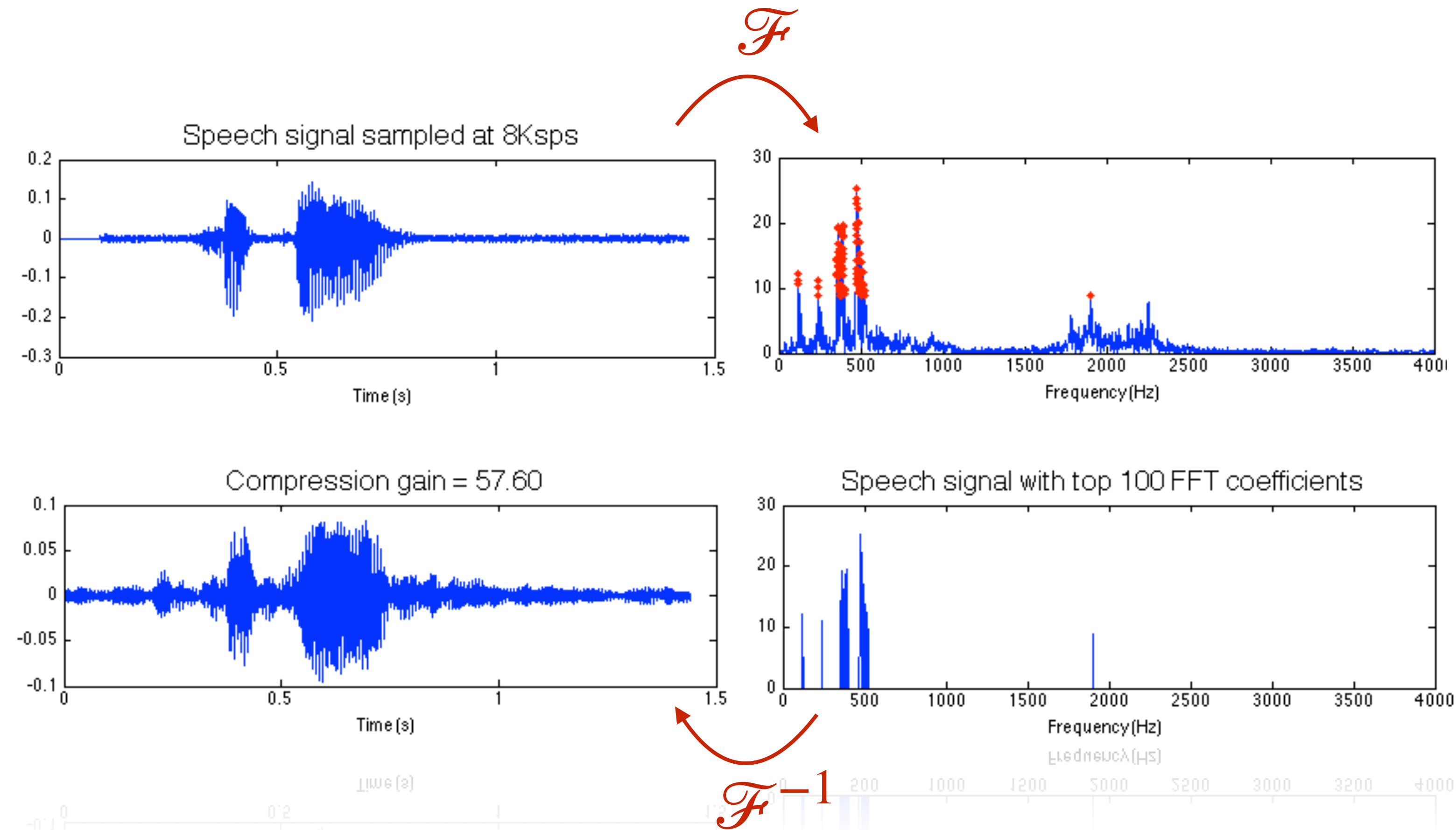
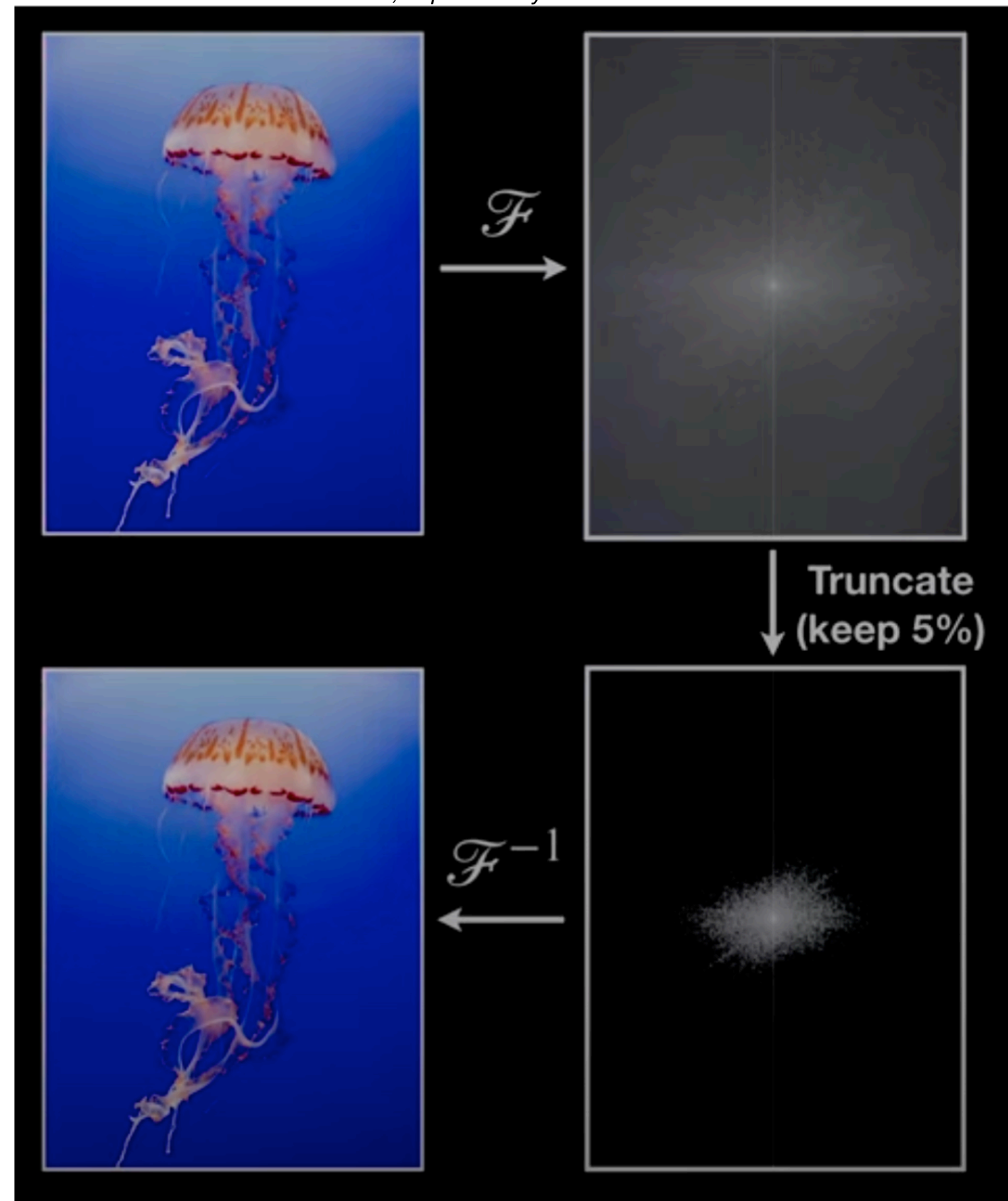
Dense signal  
(e.g. image, audio) →  $\mathbf{x}$

Fourier basis →  $\Psi$

Sparse →  $\mathbf{s}$

# Motivating Compressive Sampling

Steve Brunton, <https://www.youtube.com/watch?v=hmBTACBGWJs>



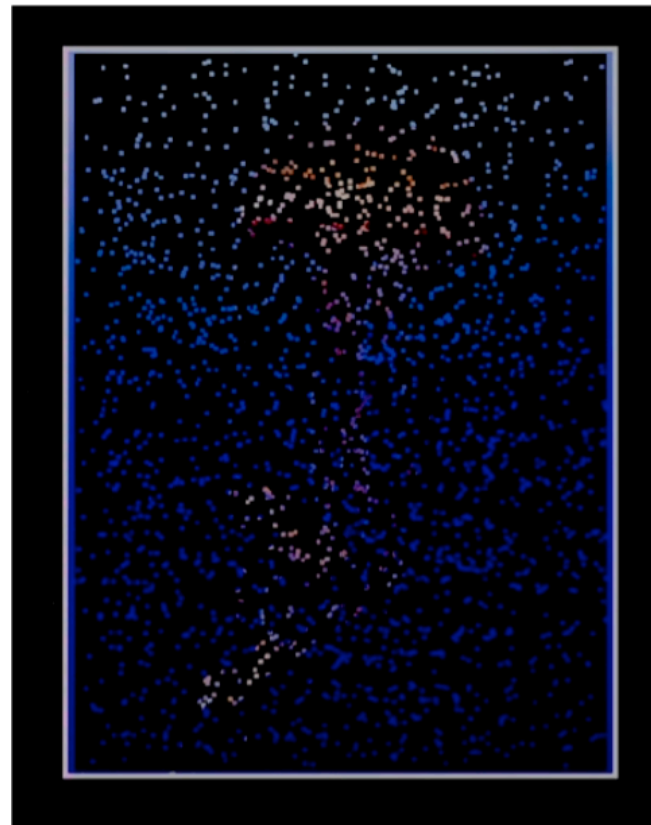
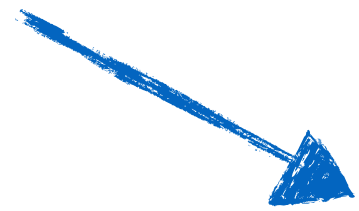
Loss compression methods exploit sparsity by retaining only top few coefficients that carry most of the signal energy.

# Motivating Compressive Sampling

---

If we throw away most of the information during compression, why do we collect it to begin with? Can we not just start with a massively downsampled signal?

Collect only  
a few samples

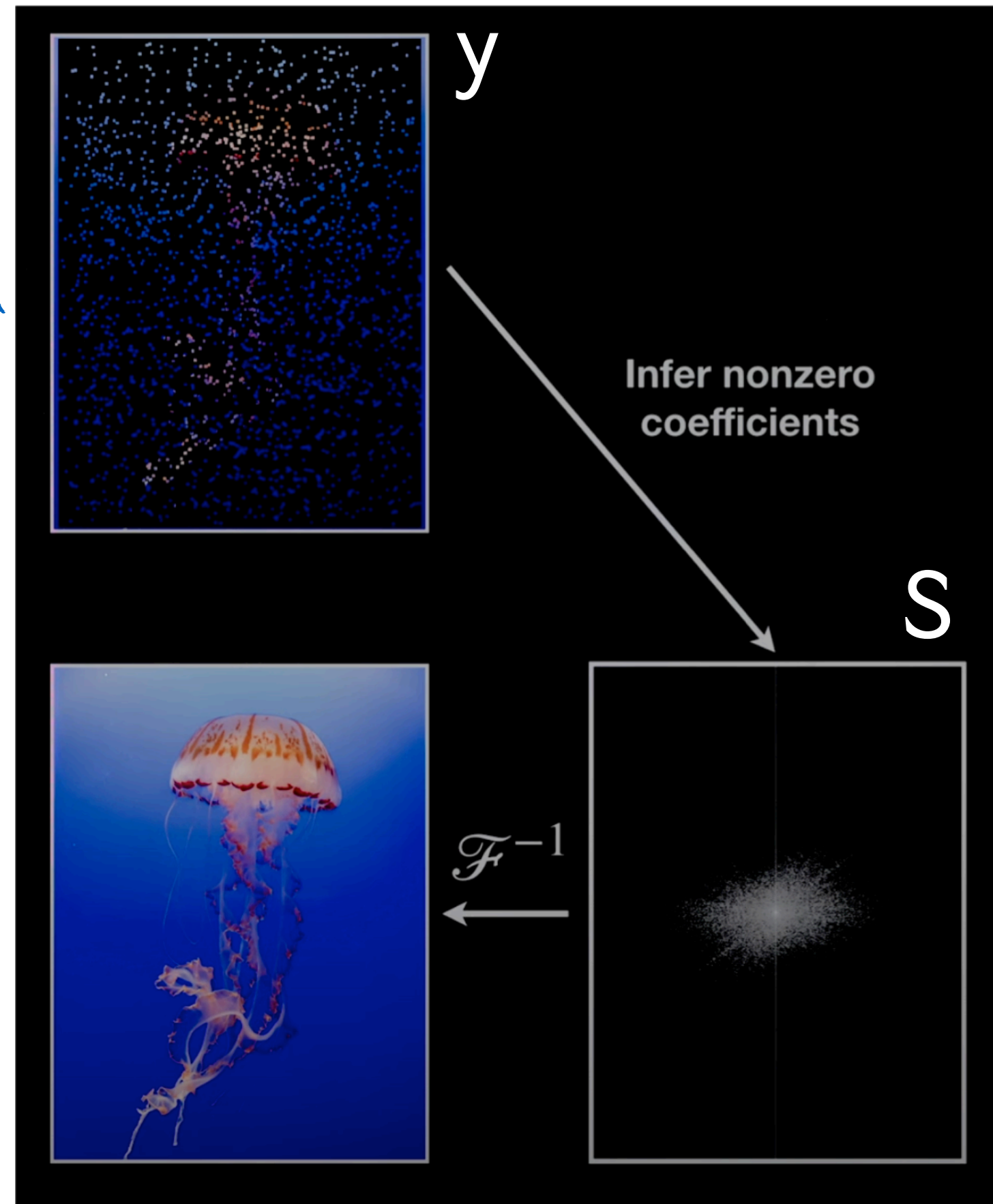


# Motivating Compressive Sampling

If we throw away most of the information during compression, why do we collect it to begin with? Can we not just start with a massively downsampled signal?

Steve Brunton, <https://www.youtube.com/watch?v=hmBTACBGWJs>

Collect only  
a few samples



Dense signal  $x = \Psi s$  Sparse

Fourier basis

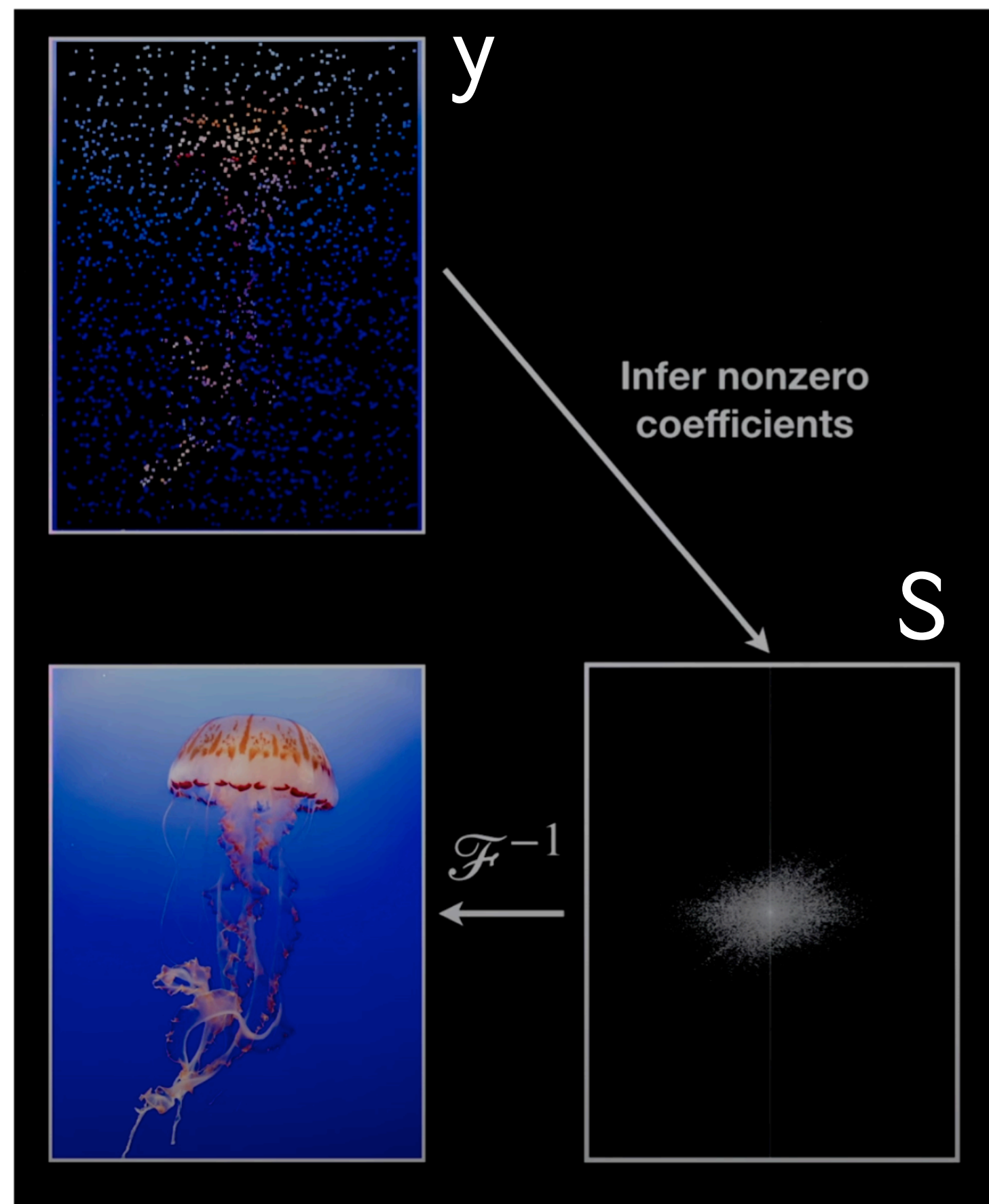
Measurements  $|y| \ll |x|$

Measurement matrix

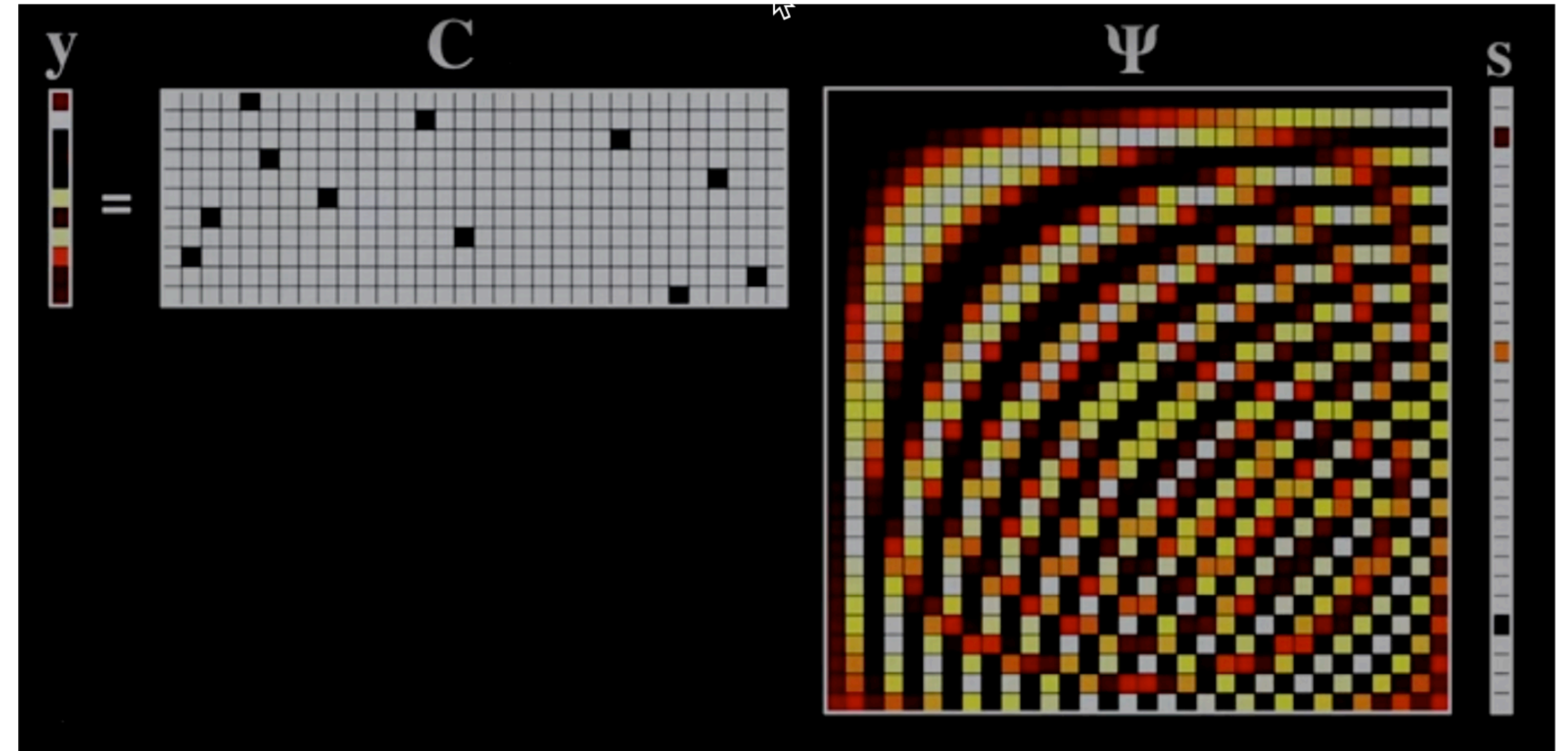
$$y = Cx$$
$$= C\Psi s$$
$$= \Theta s$$



# Inferring $s$ from $y$



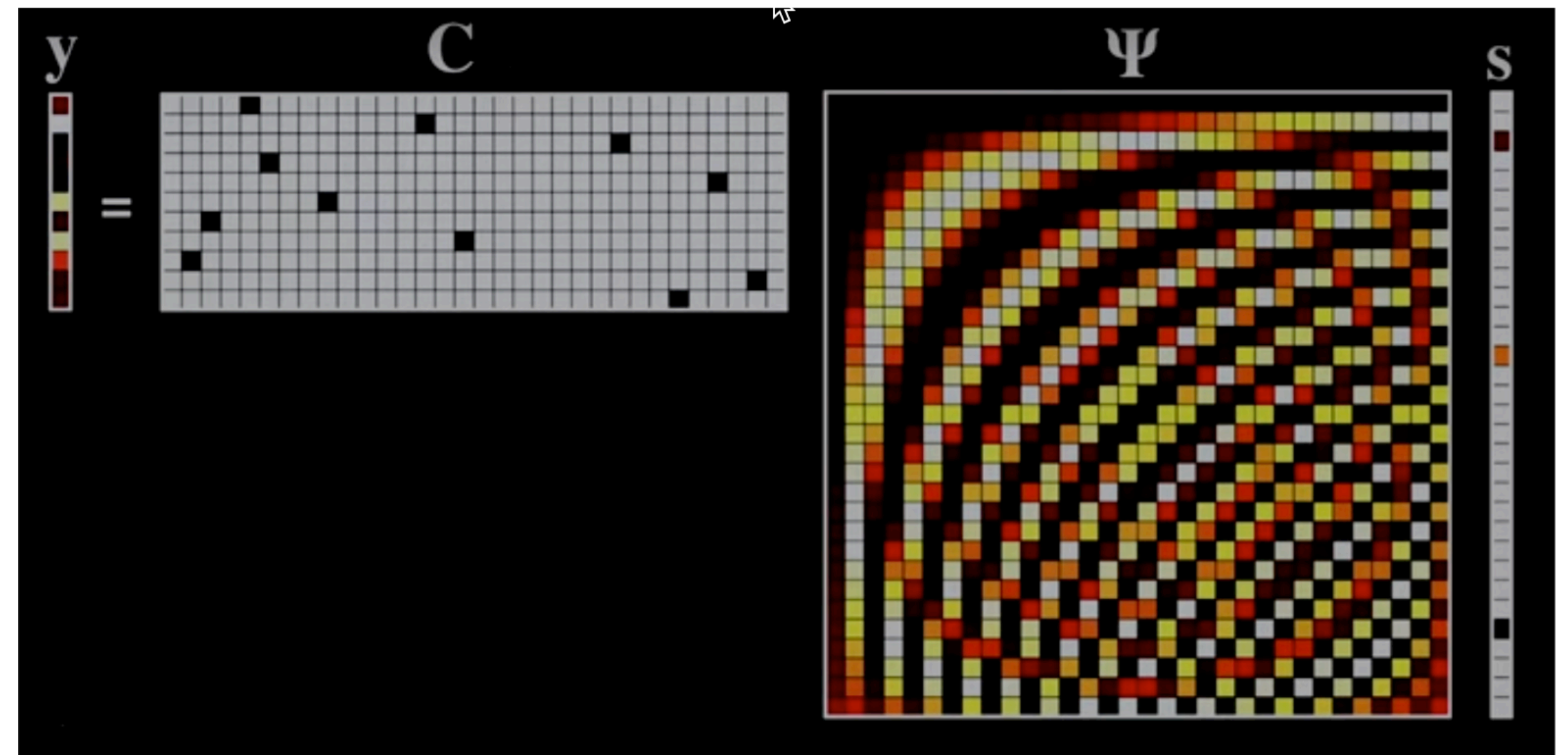
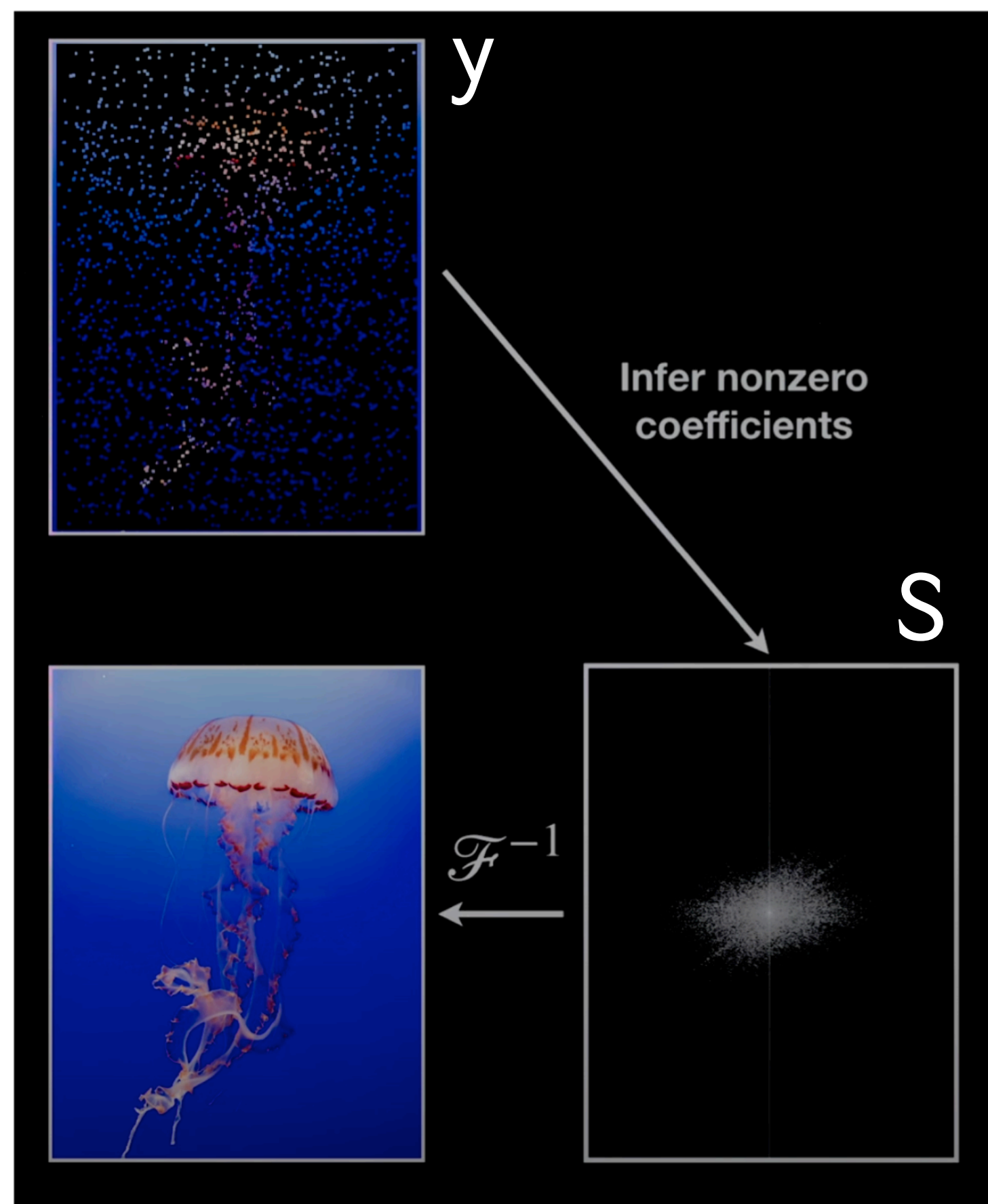
Note: once we have  $s$  we can get  $x$  via inverse-FFT



$$y = Cx = C\Psi s = \Theta s$$

**Undetermined Inverse Problem**  
Solution for  $s$  given  $y$  and  $\Theta$  is not unique.

# Adding the Sparsity Requirement to Infer $s$



$$y = Cx = C\Psi s = \Theta s$$

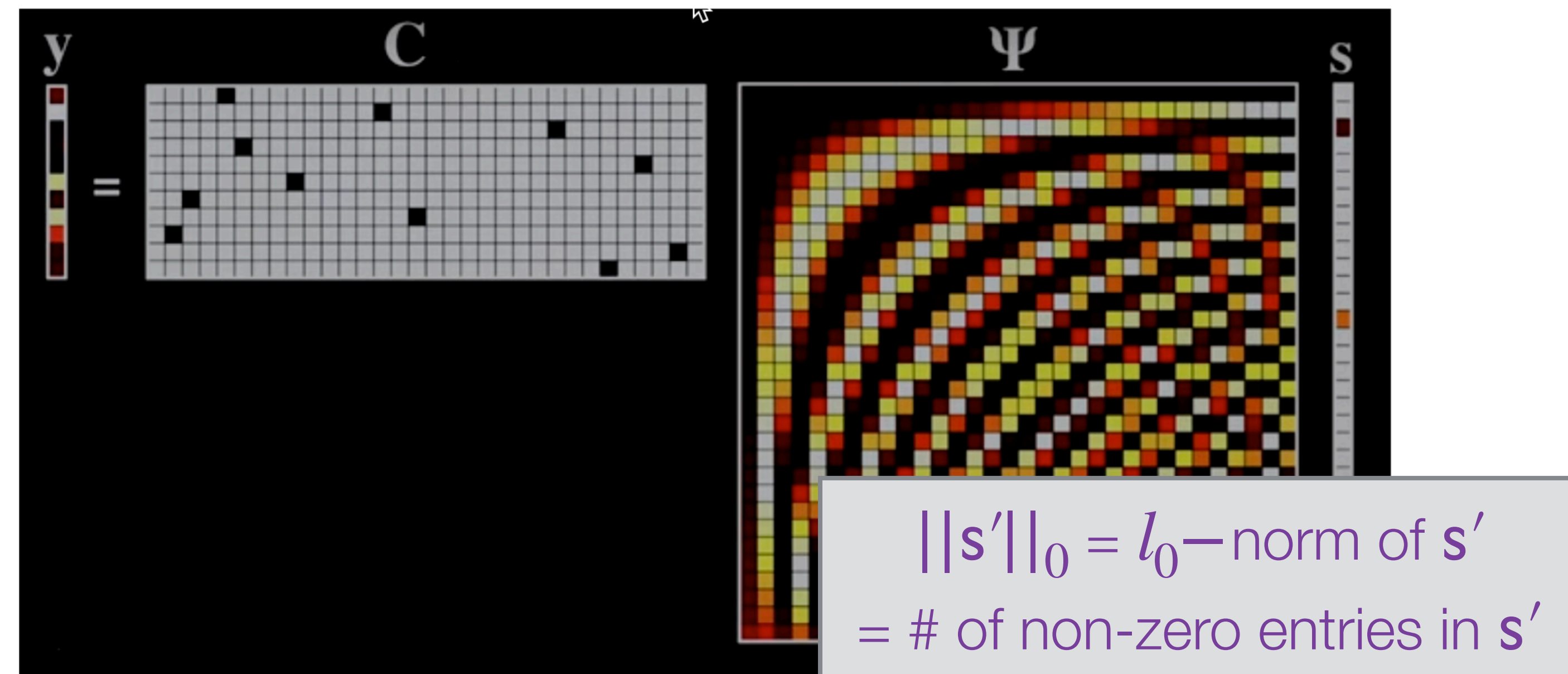
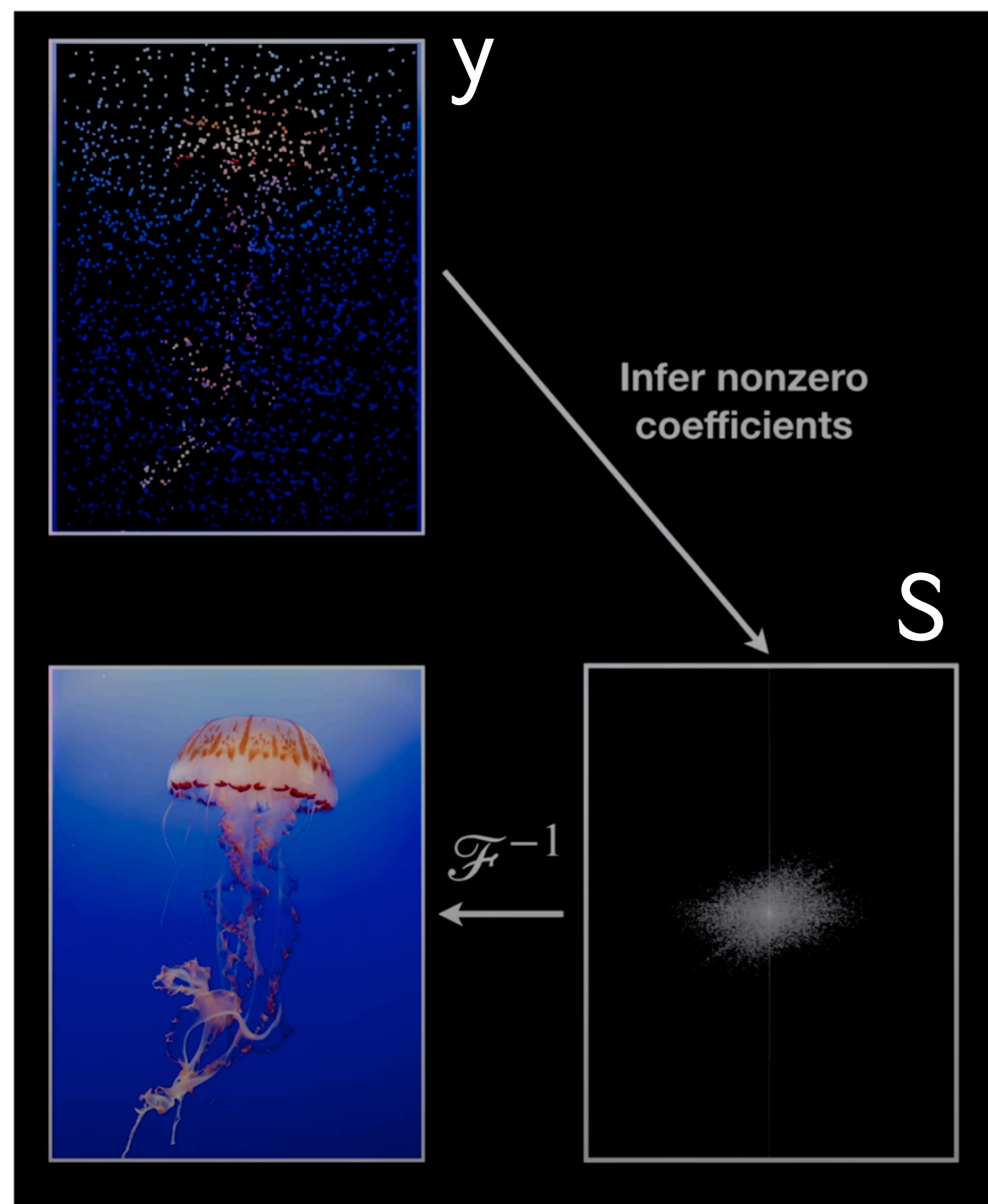
$$\hat{s} = \operatorname{argmin}_{s'} ||C\Psi s' - y||_2 + \lambda ||s'||_0$$

alternatively:

$$\hat{s} = \operatorname{argmin}_{s'} ||s'||_0, \text{ s.t. } ||C\Psi s' - y||_2 < \epsilon$$



# Adding the Sparsity Requirement to Infer $s$



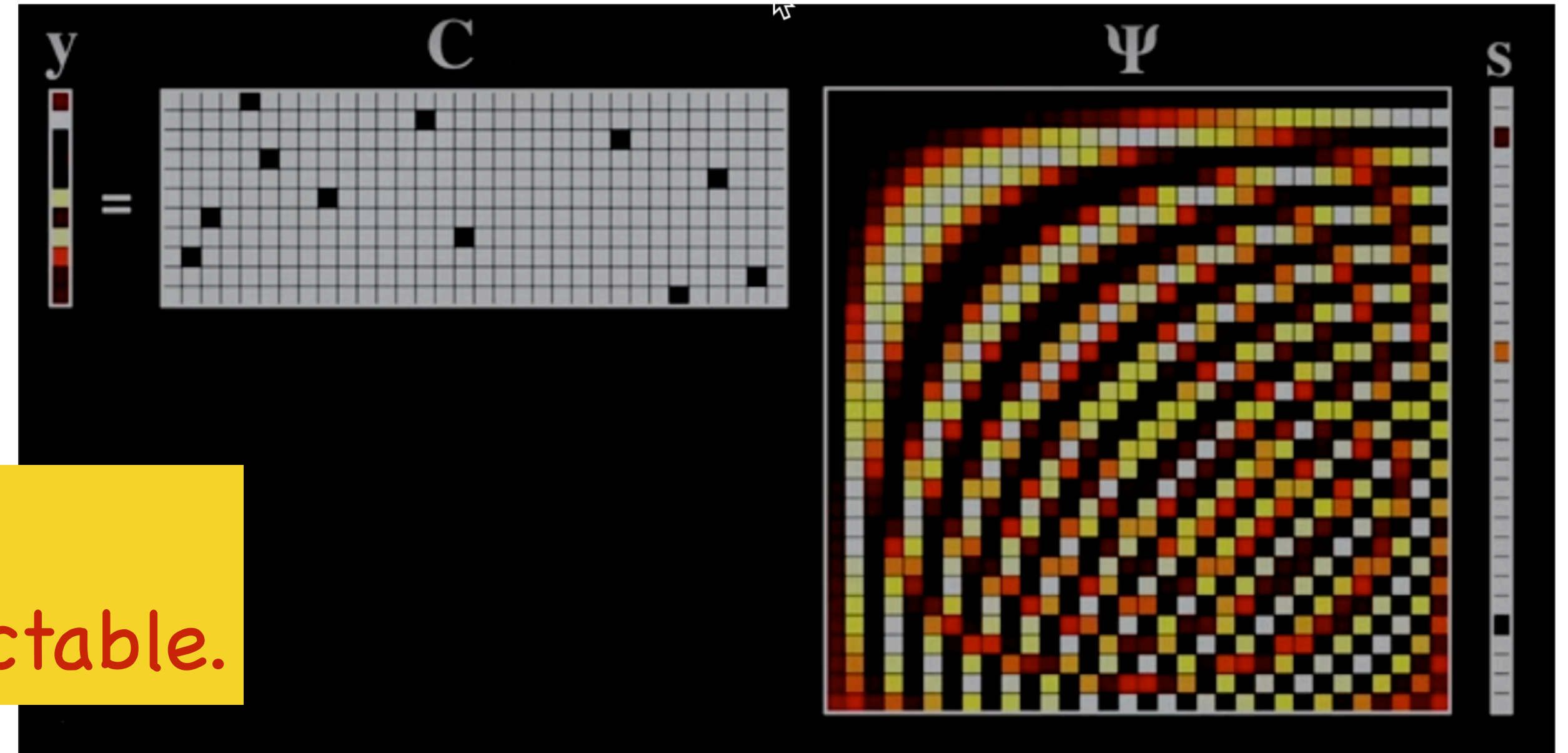
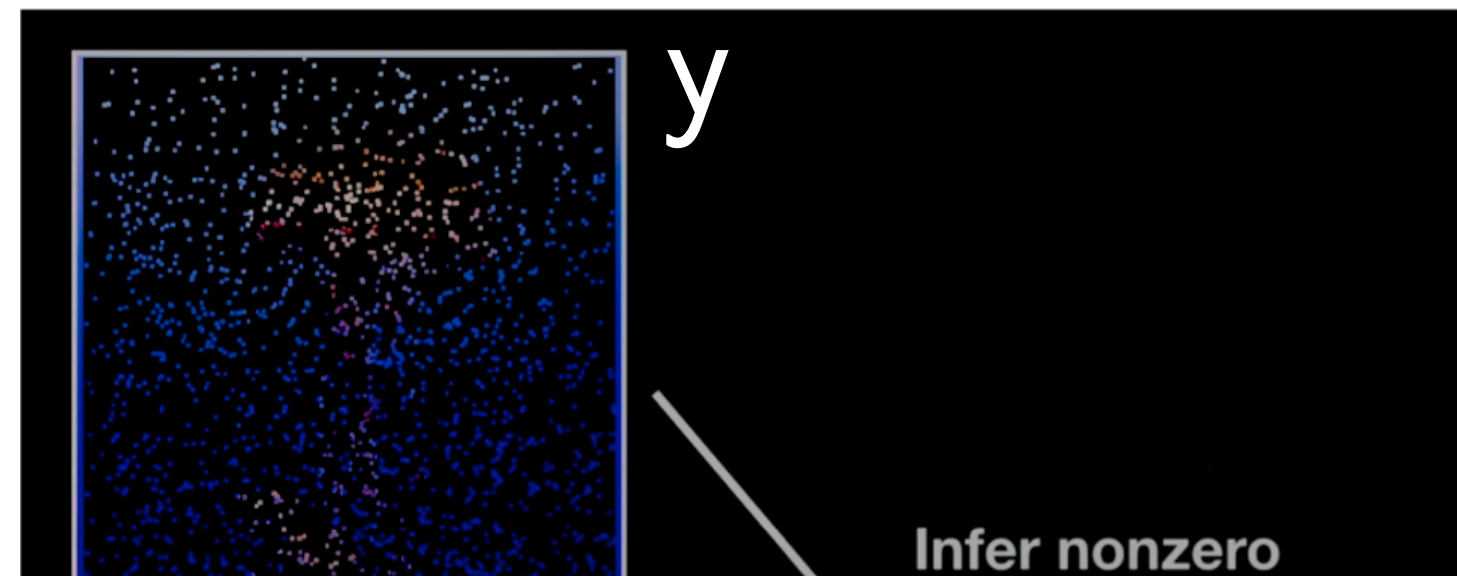
$$y = Cx = C\Psi s = \Theta s$$

$$\hat{s} = \operatorname{argmin}_{s'} \|C\Psi s' - y\|_2 + \lambda \|s'\|_0$$

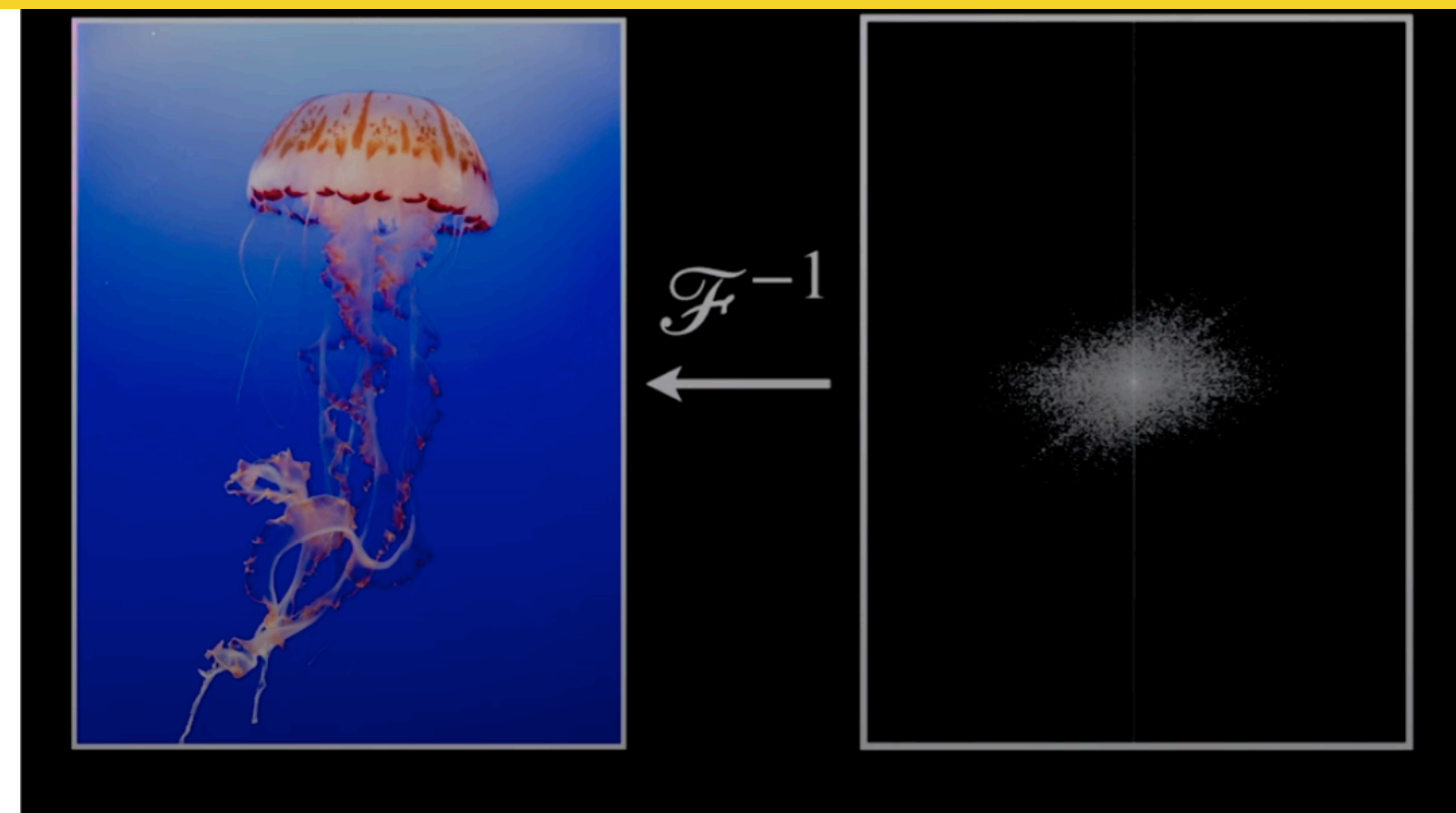
alternatively:

$$\hat{s} = \operatorname{argmin}_{s'} \|s'\|_0, \text{ s.t. } \|C\Psi s' - y\|_2 < \epsilon$$

# Adding the Sparsity Requirement to Infer $s$



**A Showstopper Issue !!!**  
This optimization problem is intractable.



$$y = Cx = C\Psi s = \Theta s$$

$$\hat{s} = \operatorname{argmin}_{s'} ||C\Psi s' - y||_2 + \lambda ||s'||_0$$

alternatively:

$$\hat{s} = \operatorname{argmin}_{s'} ||s'||_0, \text{ s.t. } ||C\Psi s' - y||_2 < \epsilon$$



# A Big Applied Math Breakthrough to the Rescue

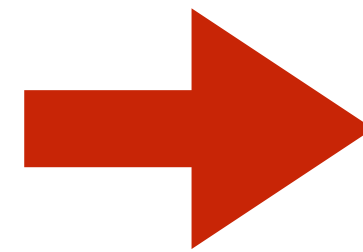
(~ 2004-2005 @ CalTech, Rice, UCLA)

$$\hat{s} = \operatorname{argmin}_{s'} ||C\Psi s' - y||_2 + \lambda ||s'||_0$$

alternatively:

$$\hat{s} = \operatorname{argmin}_{s'} ||s'||_0, \text{ s.t. } ||C\Psi s' - y||_2 < \epsilon$$

Computationally Intractable



$$\hat{s} = \operatorname{argmin}_{s'} ||C\Psi s' - y||_2 + \lambda ||s'||_1$$

alternatively:

$$\hat{s} = \operatorname{argmin}_{s'} ||s'||_1, \text{ s.t. } ||C\Psi s' - y||_2 < \epsilon$$

Computationally Efficient  
(Convex Optimization)

Gives the exact solution for  $s$   
with probability close to 1 if  
**certain conditions** are satisfied

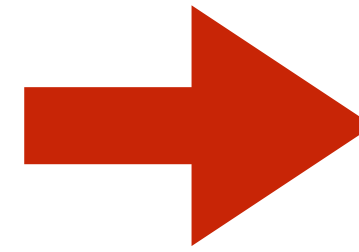
# A Big Applied Math Breakthrough to the Rescue

(~ 2004-2005 @ CalTech, Rice, UCLA)

$$\hat{s} = \operatorname{argmin}_{s'} ||C\Psi s' - y||_2 + \lambda ||s'||_0$$

alternatively:

$$\hat{s} = \operatorname{argmin}_{s'} ||s'||_0, \text{ s.t. } ||C\Psi s' - y||_2 < \epsilon$$



$$\hat{s} = \operatorname{argmin}_{s'} ||C\Psi s' - y||_2 + \lambda ||s'||_1$$

alternatively:

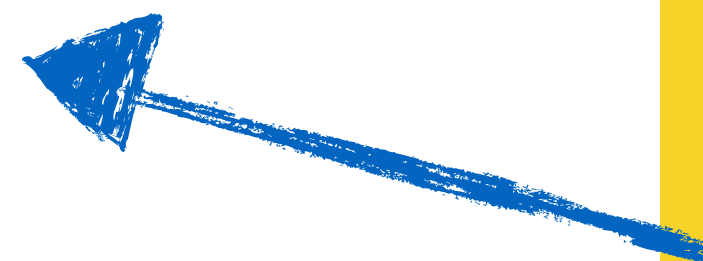
$$\hat{s} = \operatorname{argmin}_{s'} ||s'||_1, \text{ s.t. } ||C\Psi s' - y||_2 < \epsilon$$

Computationally Intractable

Computationally Efficient  
(Convex Optimization)

- ▶  $C$  should be incoherent w.r.t.  $\Psi$  (i.e. rows of  $C$  should not be too parallel to columns of  $\Psi$ )
- ▶ # of measurements  $p \sim O(k \log(n/k))$  where  $k = ||s||_0$  and  $n = |s| = |x|$

Gives the exact solution for  $s$   
with probability close to 1 if  
**certain conditions** are satisfied



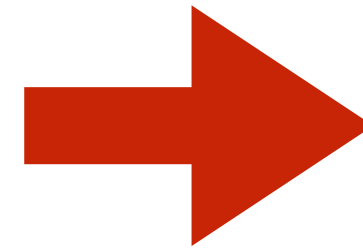
# A Big Applied Math Breakthrough to the Rescue

(~ 2004-2005 @ CalTech, Rice, UCLA)

$$\hat{s} = \operatorname{argmin}_{s'} ||C\Psi s' - y||_2 + \lambda ||s'||_0$$

alternatively:

$$\hat{s} = \operatorname{argmin}_{s'} ||s'||_0, \text{ s.t. } ||C\Psi s' - y||_2 < \epsilon$$



$$\hat{s} = \operatorname{argmin}_{s'} ||C\Psi s' - y||_2 + \lambda ||s'||_1$$

alternatively:

$$\hat{s} = \operatorname{argmin}_{s'} ||s'||_1, \text{ s.t. } ||C\Psi s' - y||_2 < \epsilon$$

Computationally Intractable

Computationally Efficient  
(Convex Optimization)

- ▶  $C$  should be incoherent w.r.t.  $\Psi$  (i.e. rows of  $C$  should not be too parallel to columns of  $\Psi$ )
- ▶ # of measurements  $p \sim O(k \log(n/k))$  where  $k = ||s||_0$  and  $n = |s| = |x|$

Gives the exact solution for  $s$   
with probability close to 1 if  
**certain conditions** are satisfied

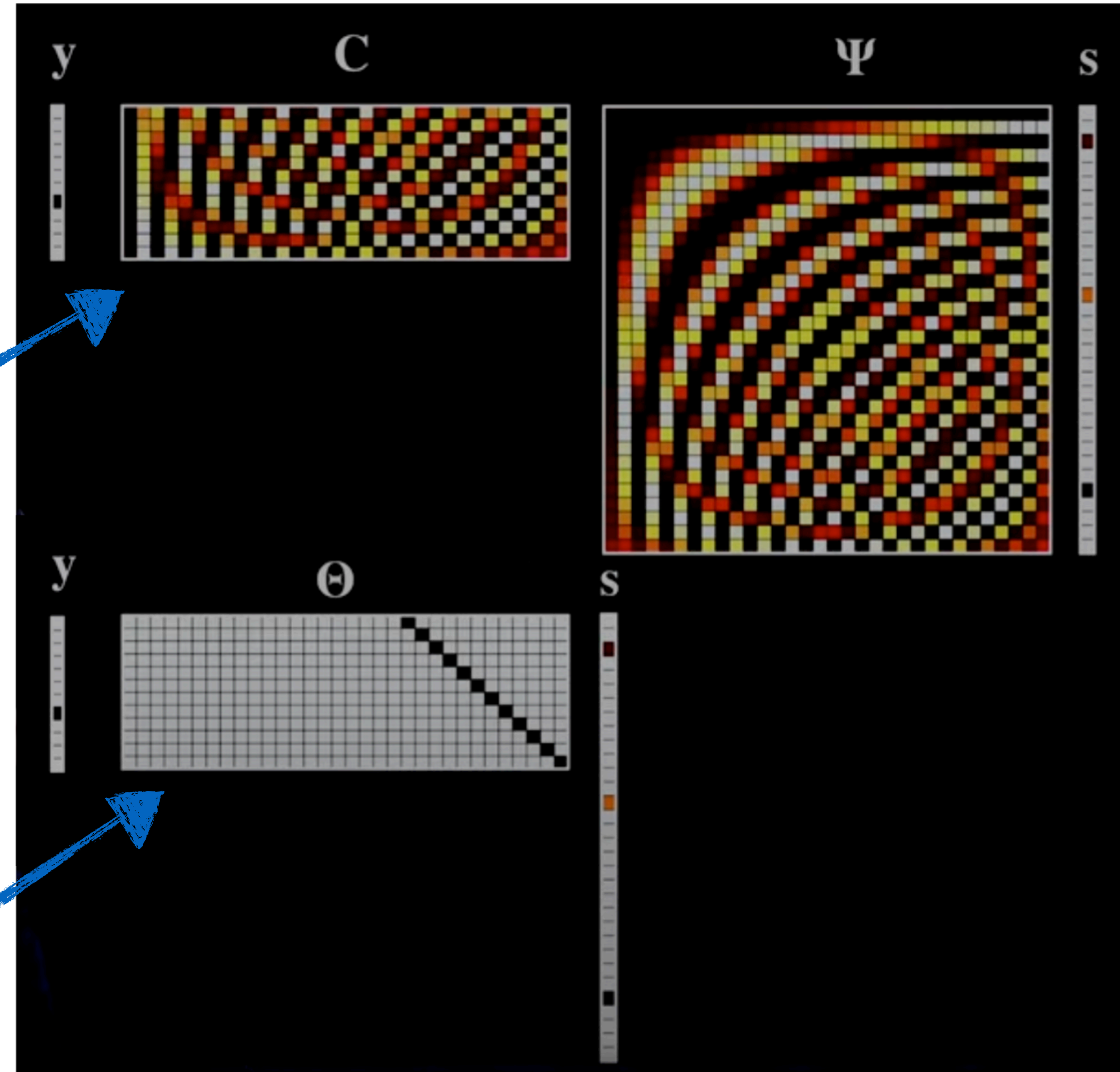
Restricted Isometry Property (RIP)

i.e.  $C\Psi$  acts like a unitary matrix on sparse vector  $s$

# Example of Bad Measurement Matrix $C$

$C$  has rows that are the same as a subset of columns of  $\Psi$

$\Theta$  only picks up information only from a part of  $s$

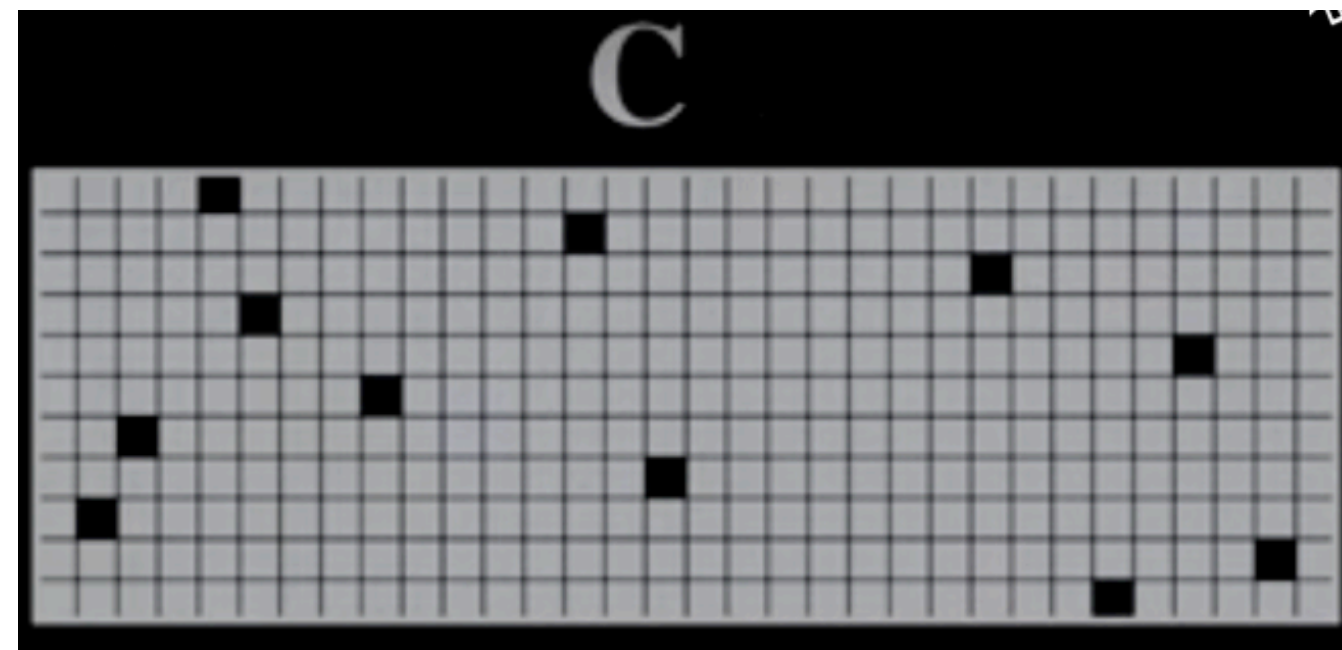




# Measurement Matrix $C$ in Practice

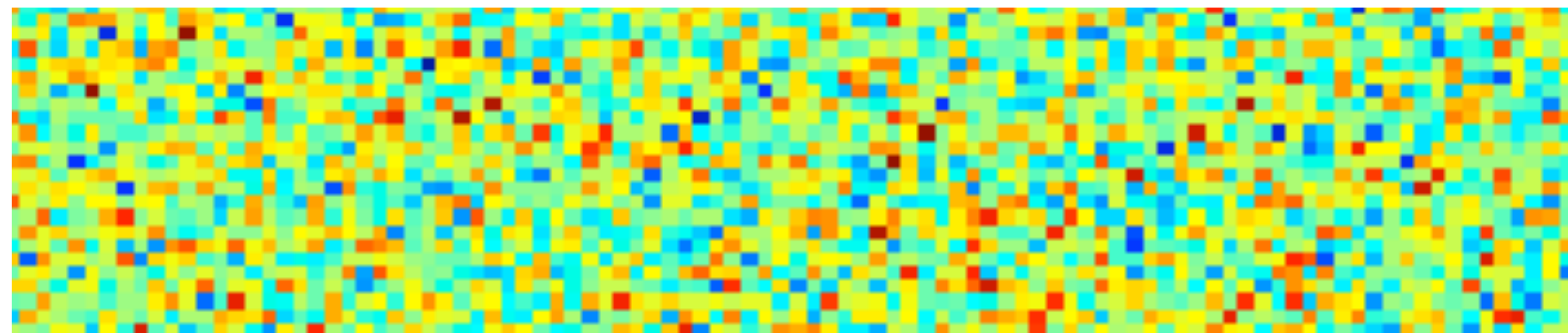
---

- **Random 0-1 mask matrix:** each element in  $y$  is a random sample picked from  $x$ 
  - ▶ rows of  $C$  have exactly one 1 to indicate the sample picked by that measurement
  - ▶ must be causal in case of real-time sampling of a time-series



Steve Brunton, <https://www.youtube.com/watch?v=hmBTACBGWJs>

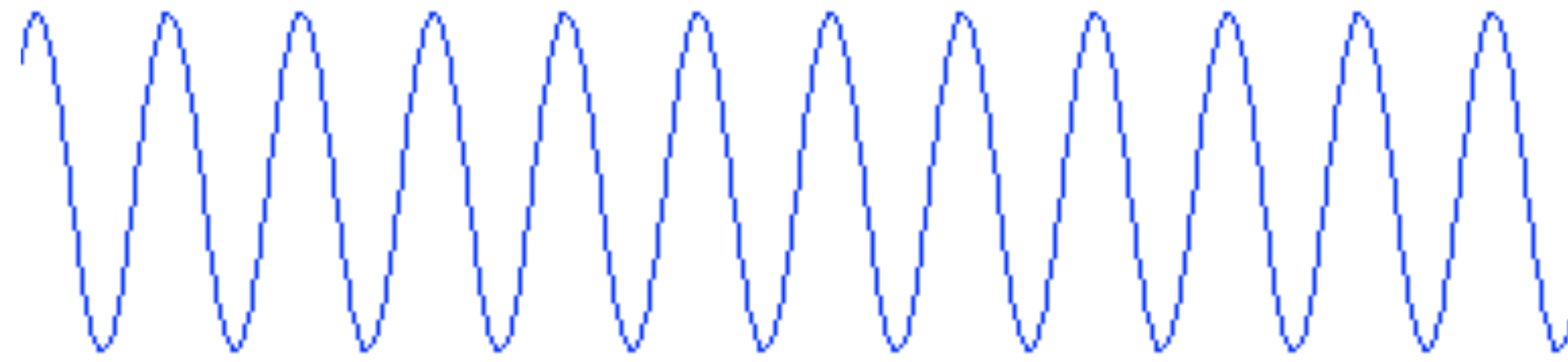
- **Random real-numbers:** each element in  $y$  is an incoherent random projection of  $x$ 
  - ▶ commonly used: Gaussian, Bernoulli



# Causal Random Sampling of a Time Series

---

**x**

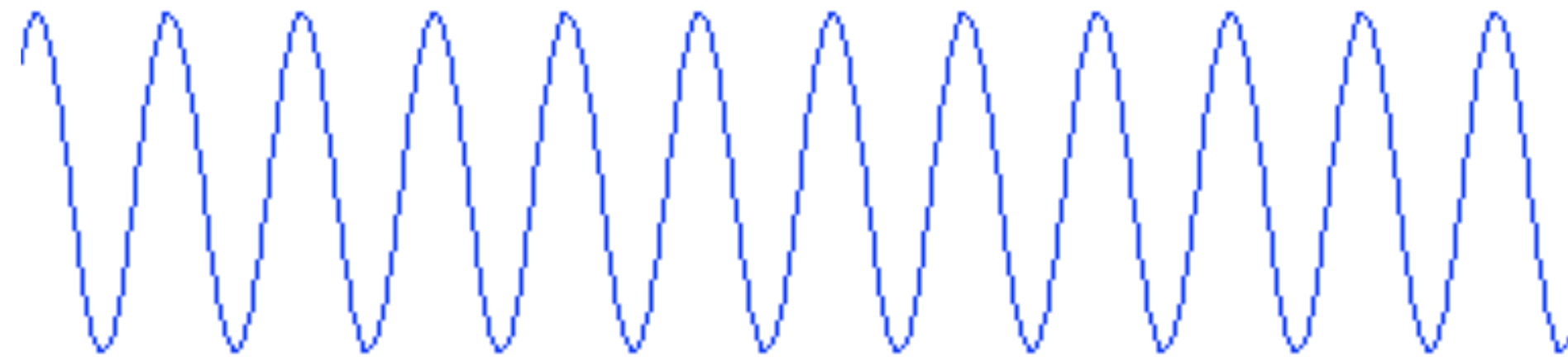


**Signal to be  
acquired**

# Causal Random Sampling of a Time Series

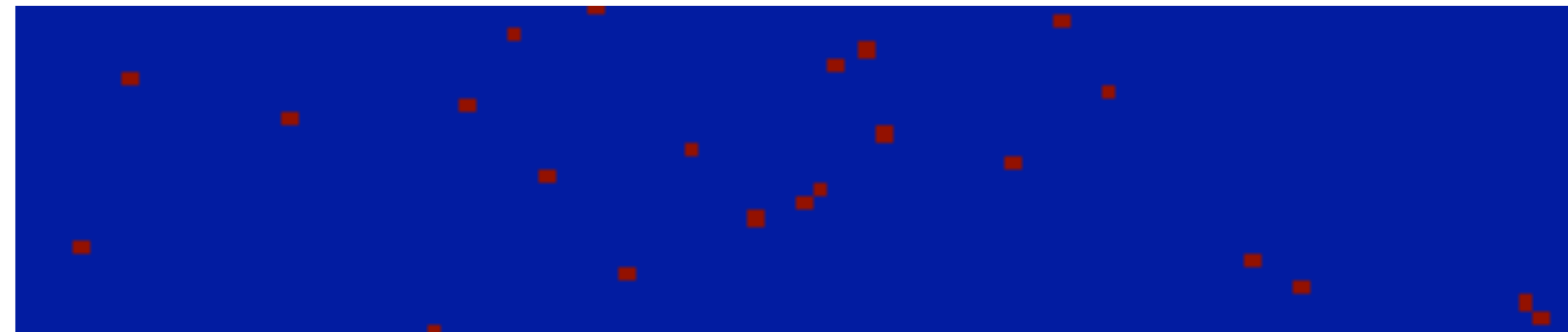
---

**X**



**Signal to be  
acquired**

**C**

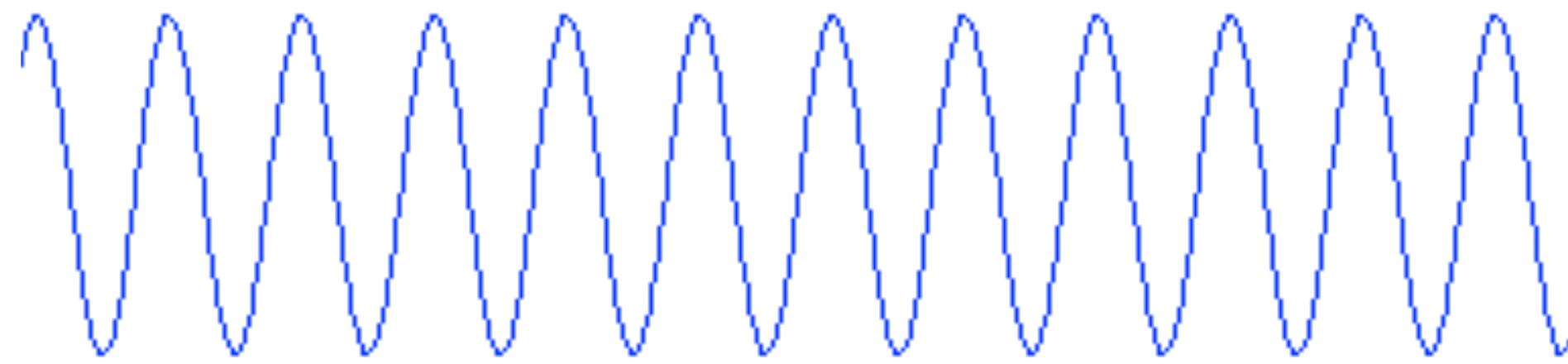


**100 x 100 identity matrix with  
75 rows uniformly randomly  
selected and thrown out**

# Causal Random Sampling of a Time Series

---

$x$



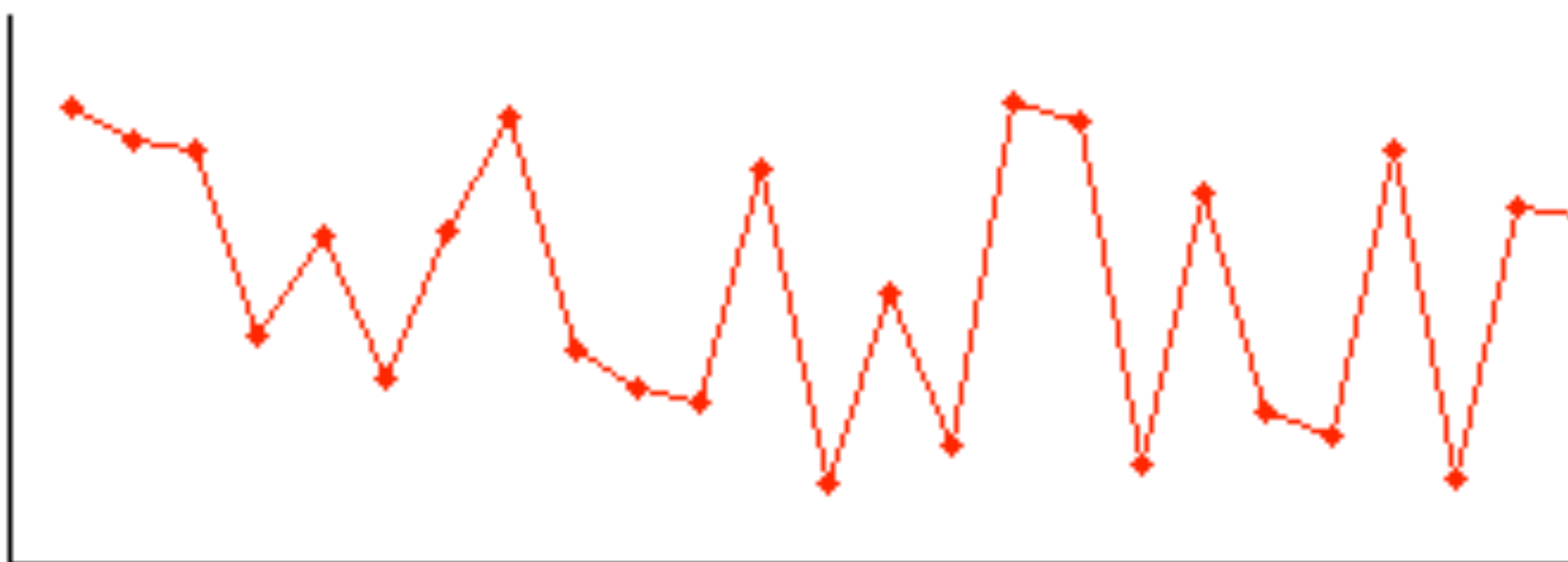
Signal to be  
acquired

$C$



100 x 100 identity matrix with  
75 rows uniformly randomly  
selected and thrown out

$y$



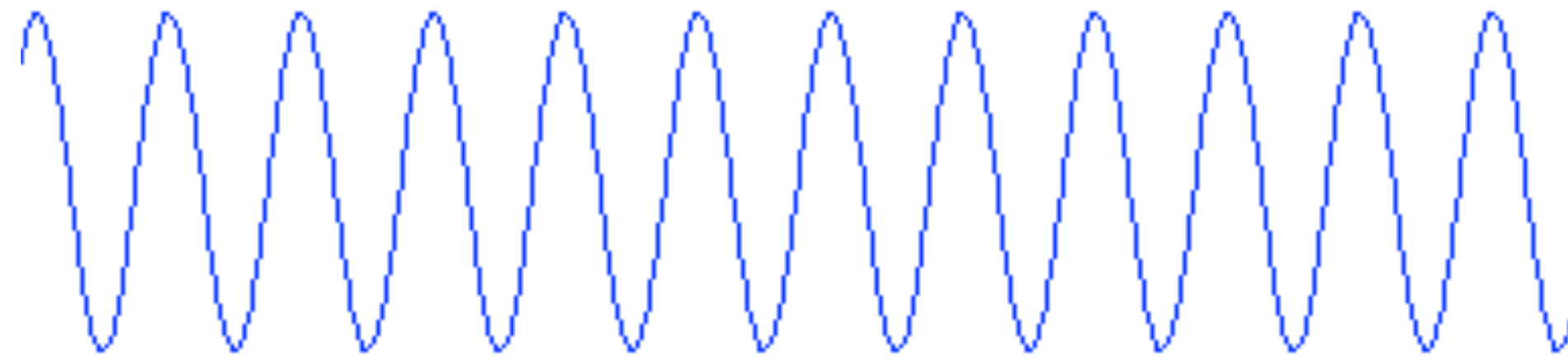
Measurements  
actually acquired



# Taking Incoherent Projections of a Time Series - Gaussian

---

**x**

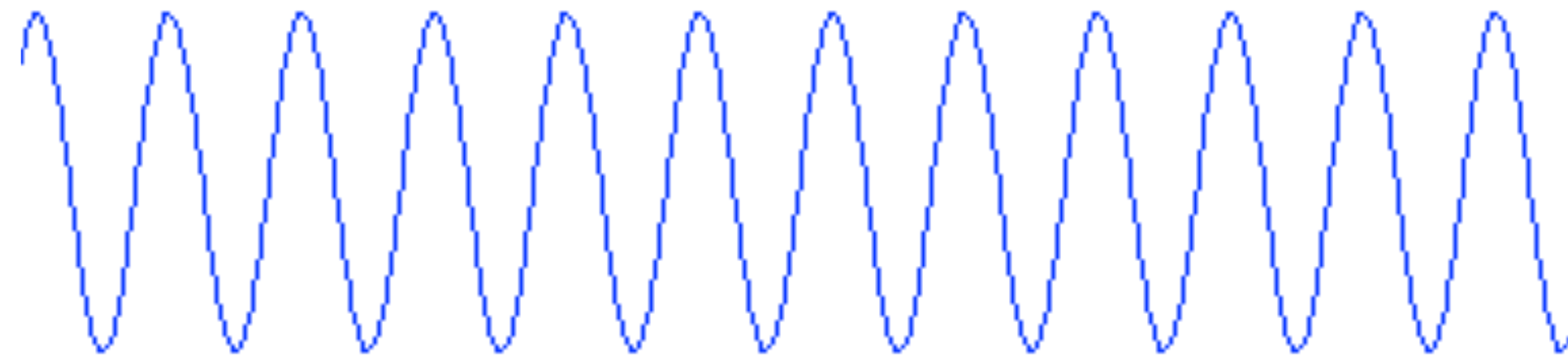


**Signal to be  
acquired**

# Taking Incoherent Projections of a Time Series - Gaussian

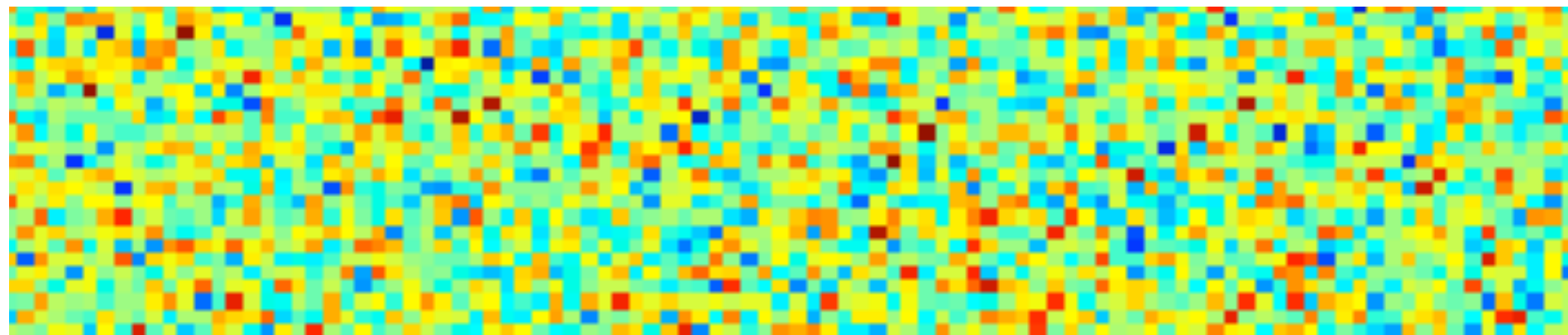
---

**x**



Signal to be  
acquired

**C**

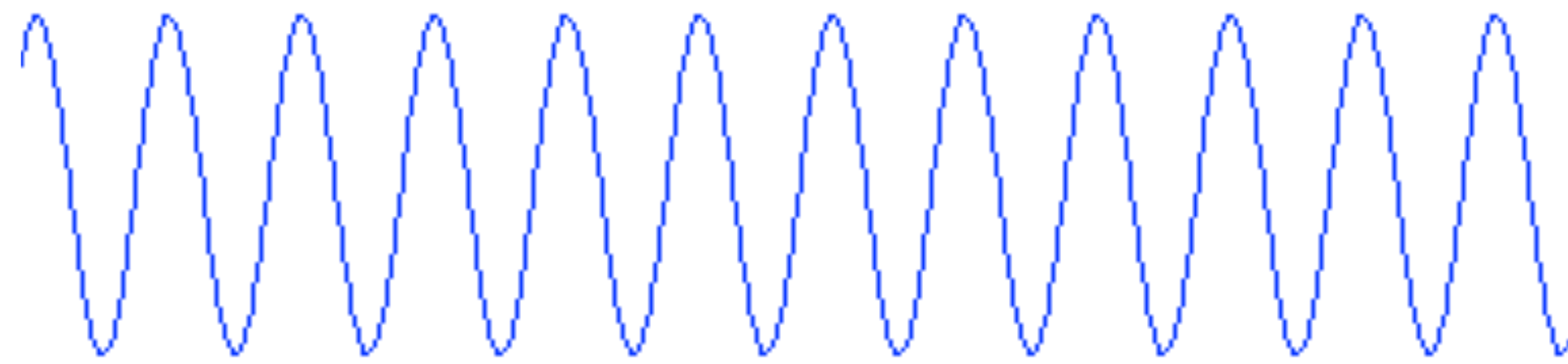


25 x 100 matrix of Gaussian  $\mathcal{N}(0, \frac{1}{n})$   
distributed random numbers

# Taking Incoherent Projections of a Time Series - Gaussian

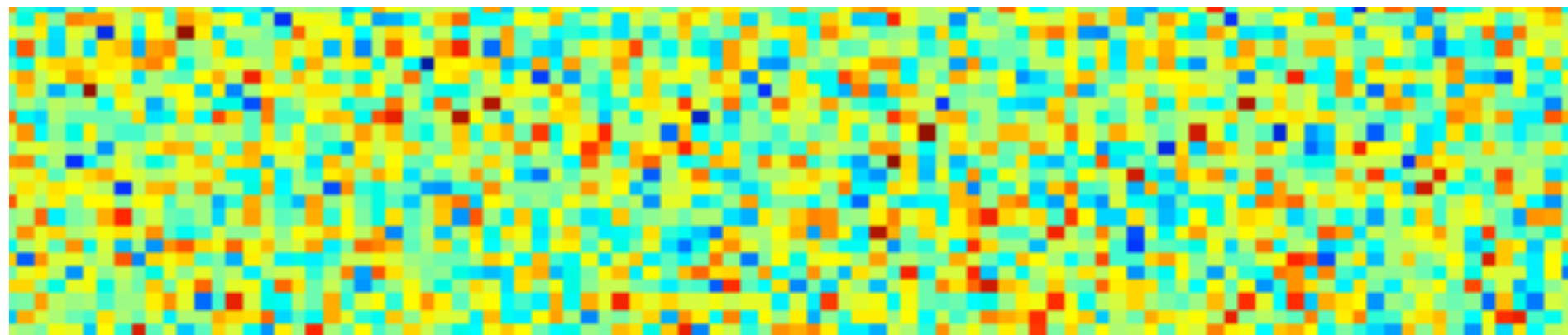
---

**x**



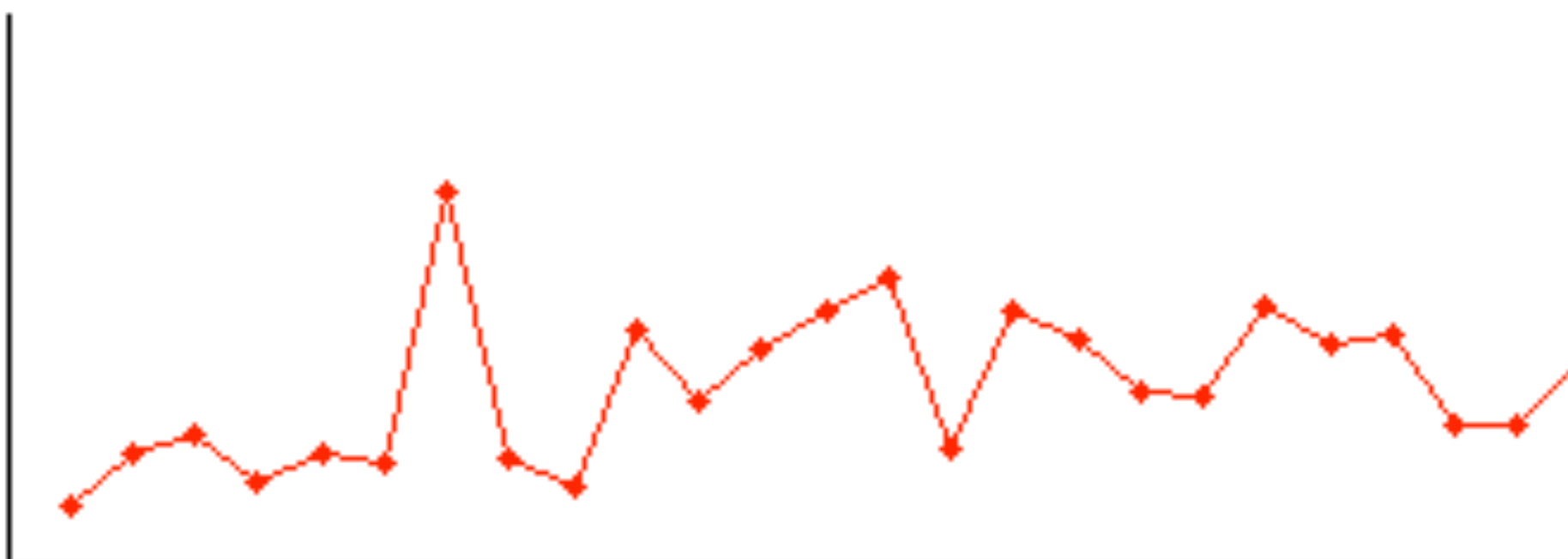
Signal to be  
acquired

**C**



25 x 100 matrix of Gaussian  $\mathcal{N}(0, \frac{1}{n})$   
distributed random numbers

**y**

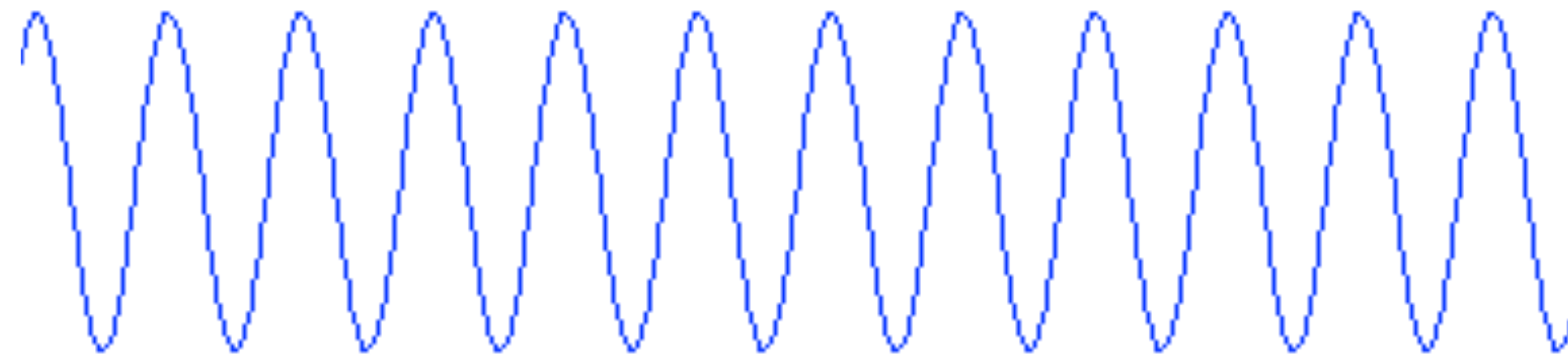


Measurements  
actually acquired

# Taking Incoherent Projections of a Time Series - Bernoulli

---

**x**

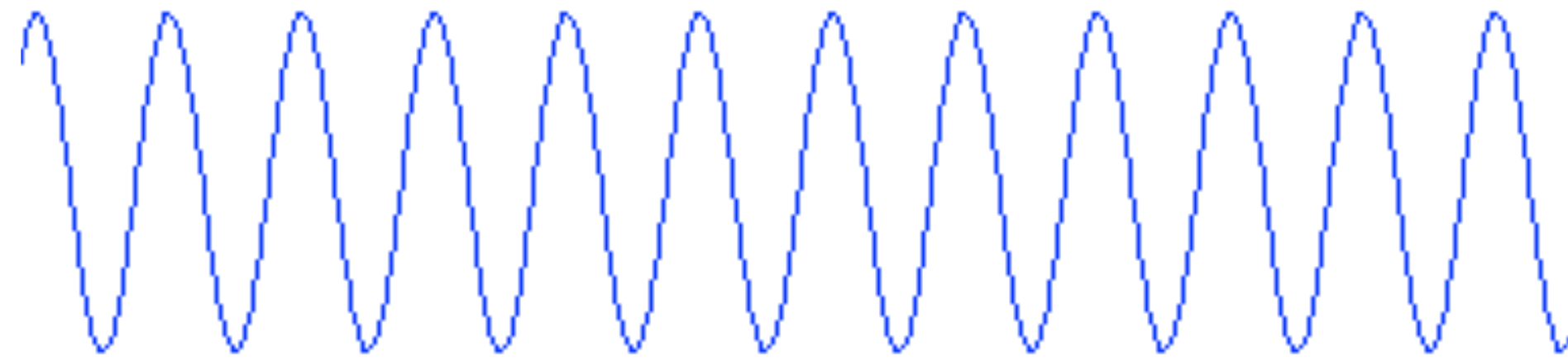


**Signal to be  
acquired**

# Taking Incoherent Projections of a Time Series - Bernoulli

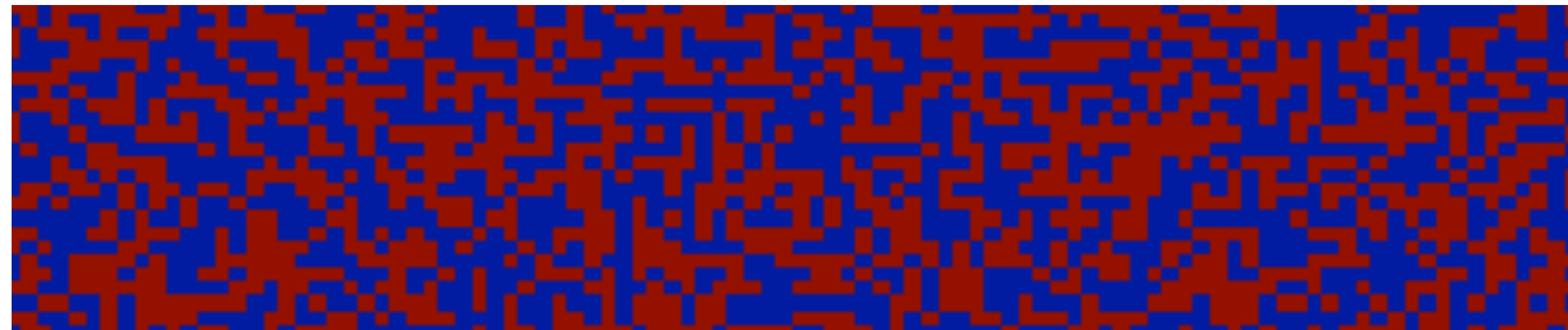
---

**x**



Signal to be  
acquired

**C**

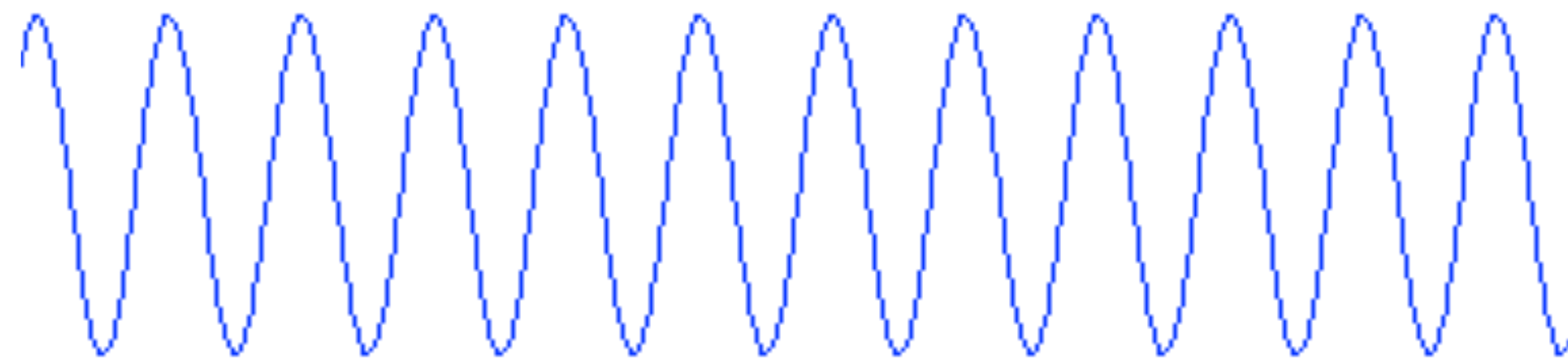


25 x 100 Bernoulli distributed ( $p = 0.5$ )  
random numbers  $\in \left\{ -\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}} \right\}$

# Taking Incoherent Projections of a Time Series - Bernoulli

---

**x**



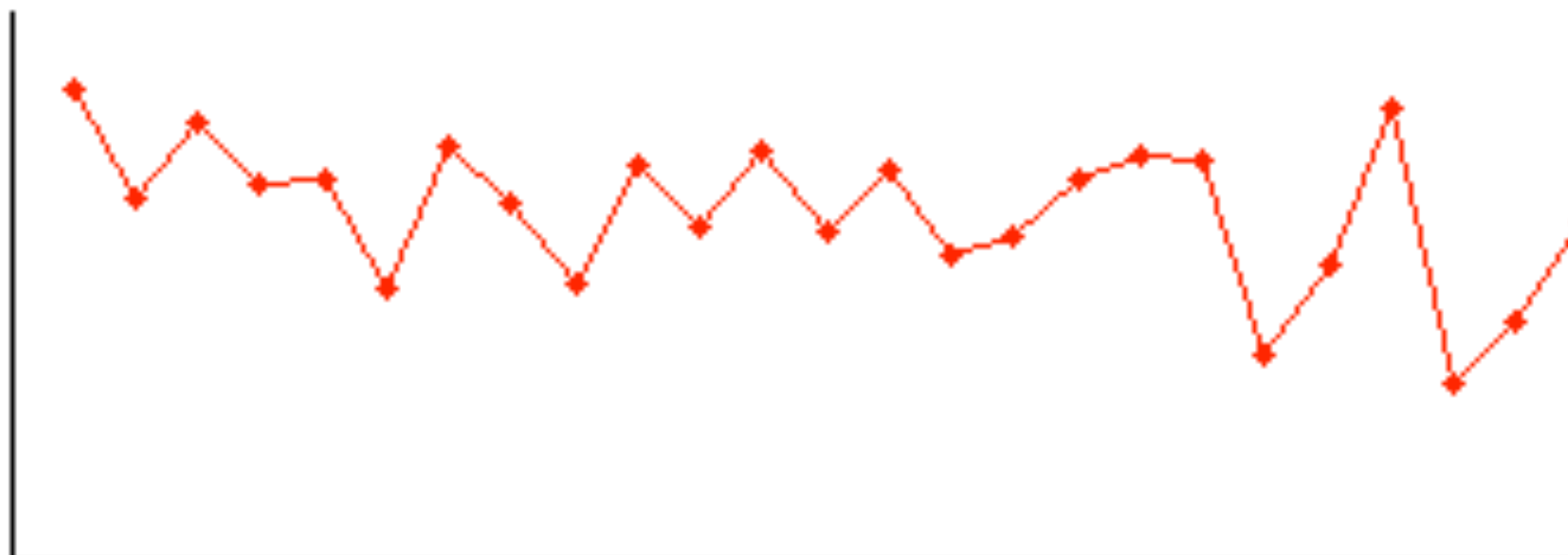
Signal to be  
acquired

**C**



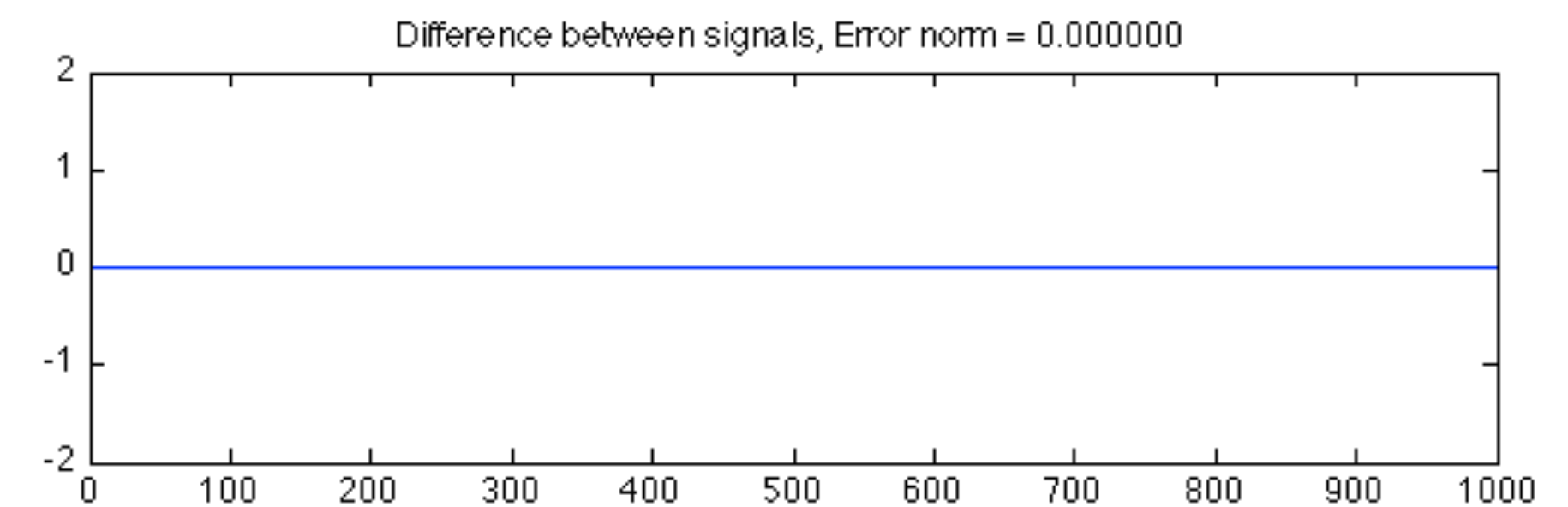
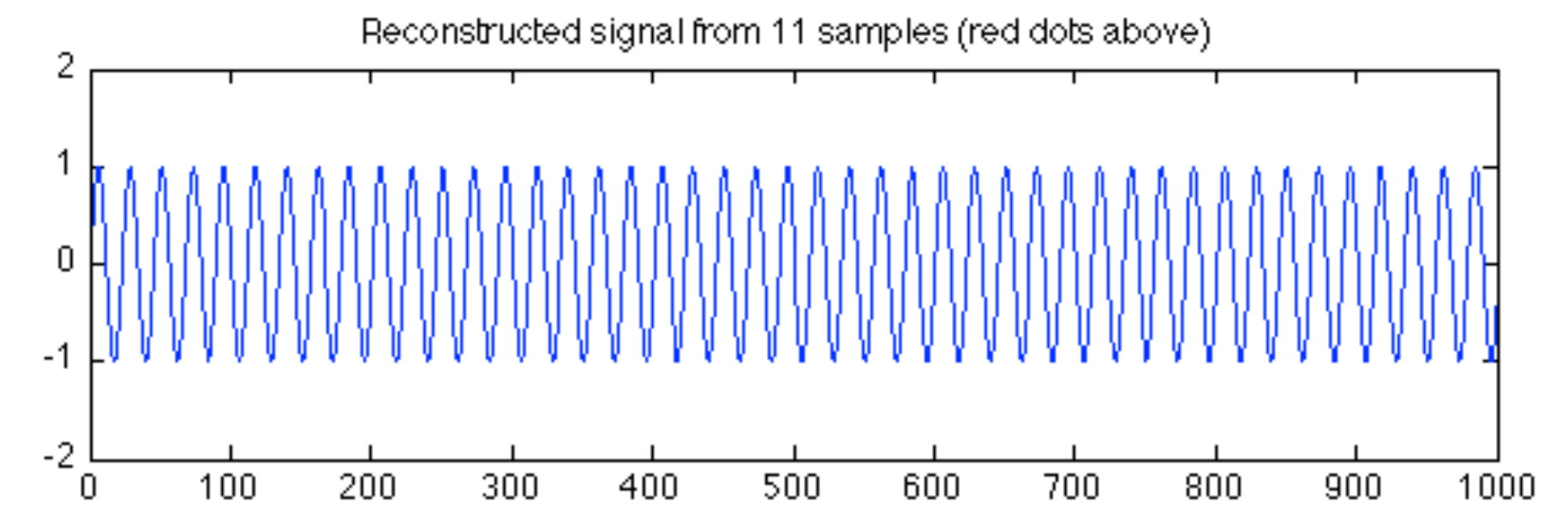
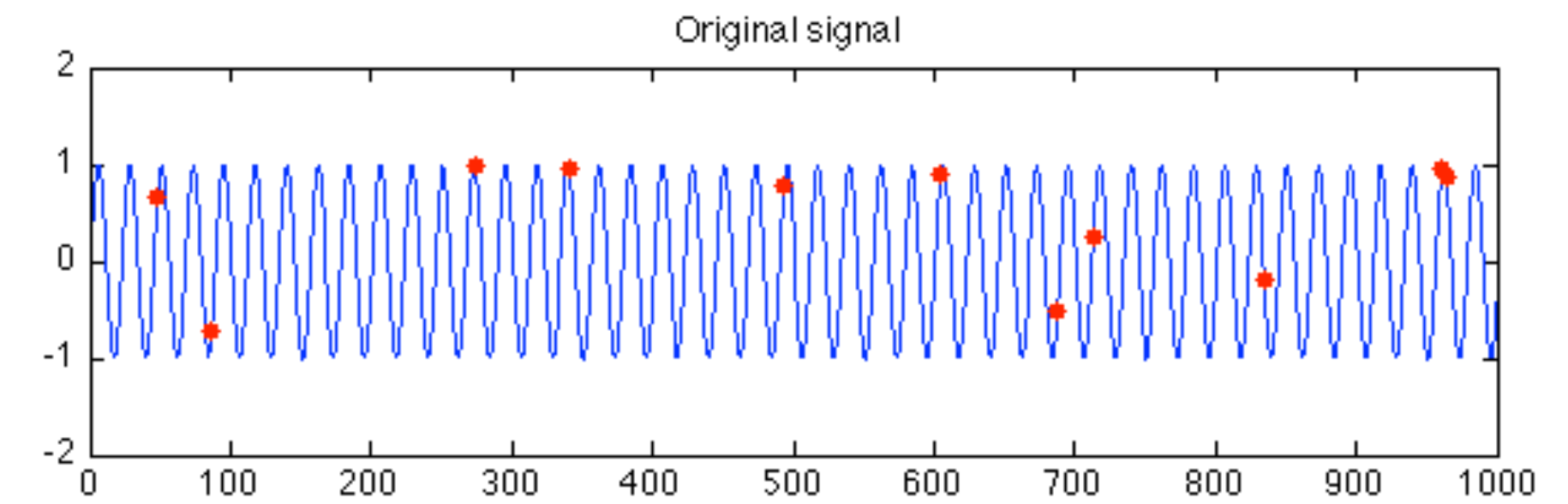
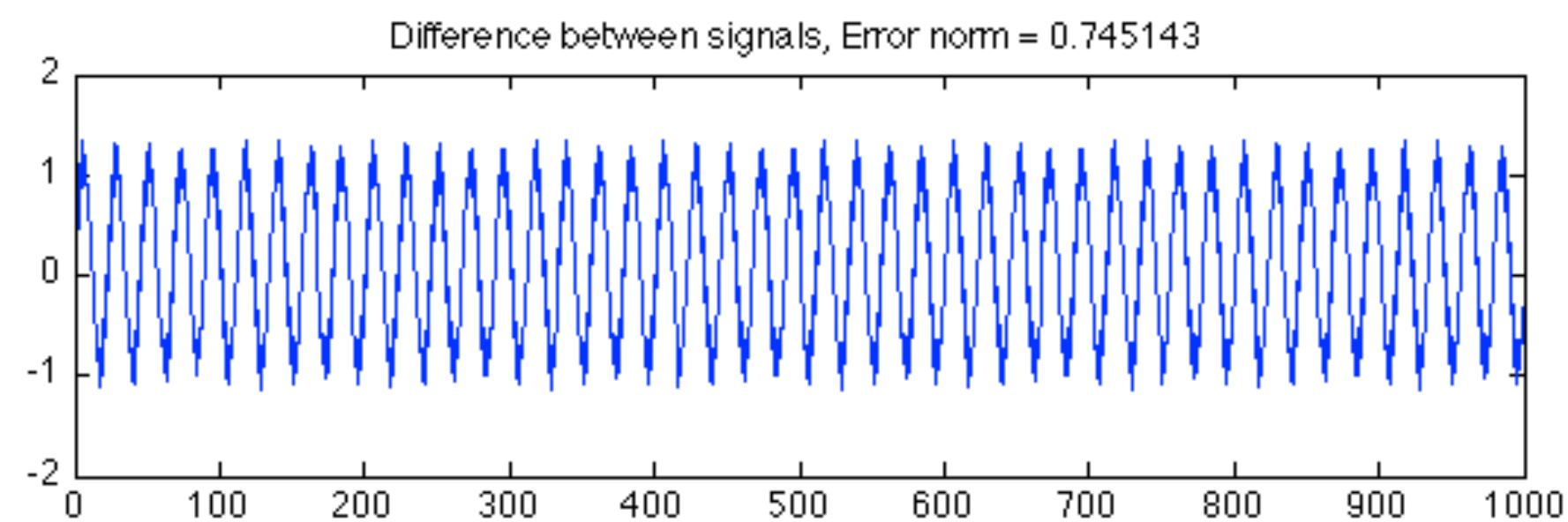
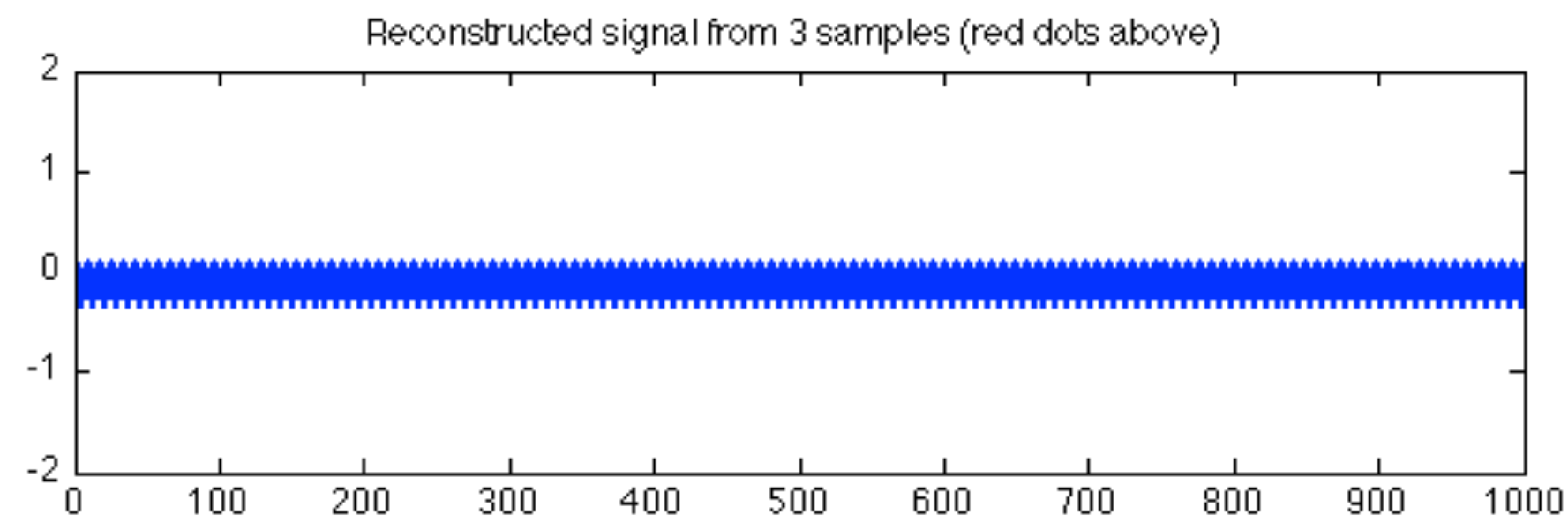
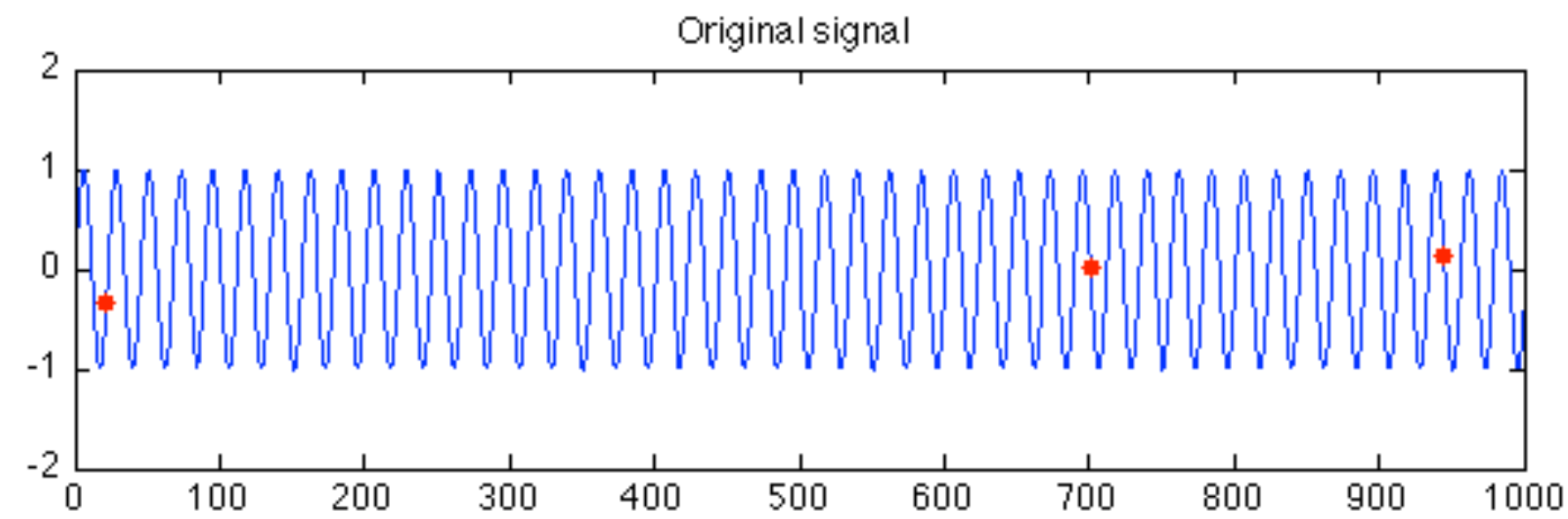
25 x 100 Bernoulli distributed ( $p = 0.5$ )  
random numbers  $\in \left\{ -\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}} \right\}$

**y**



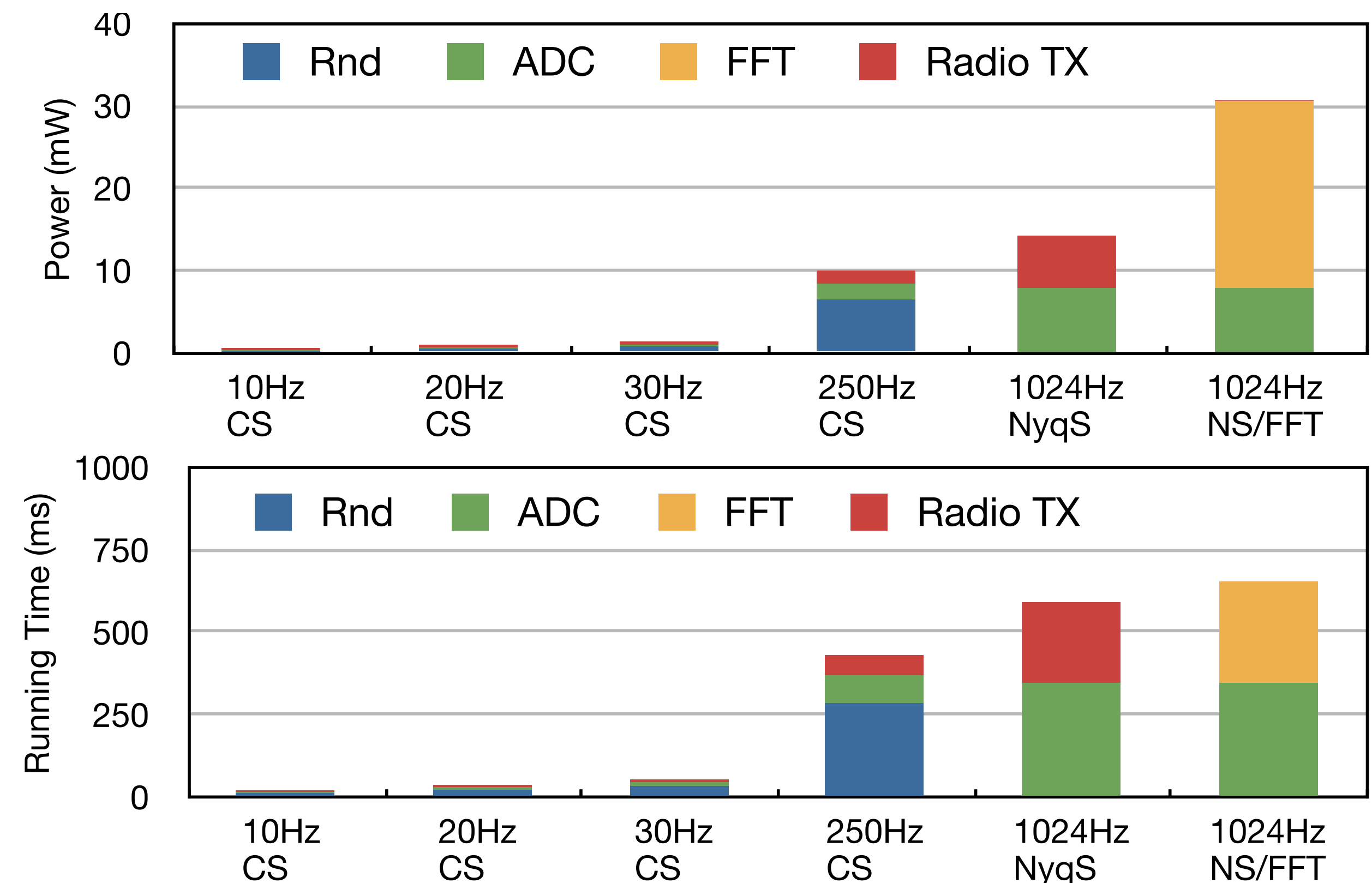
Measurements  
actually acquired

# Compressed Sensing Example



# Compressive Sampling to Save Energy on Edge Devices

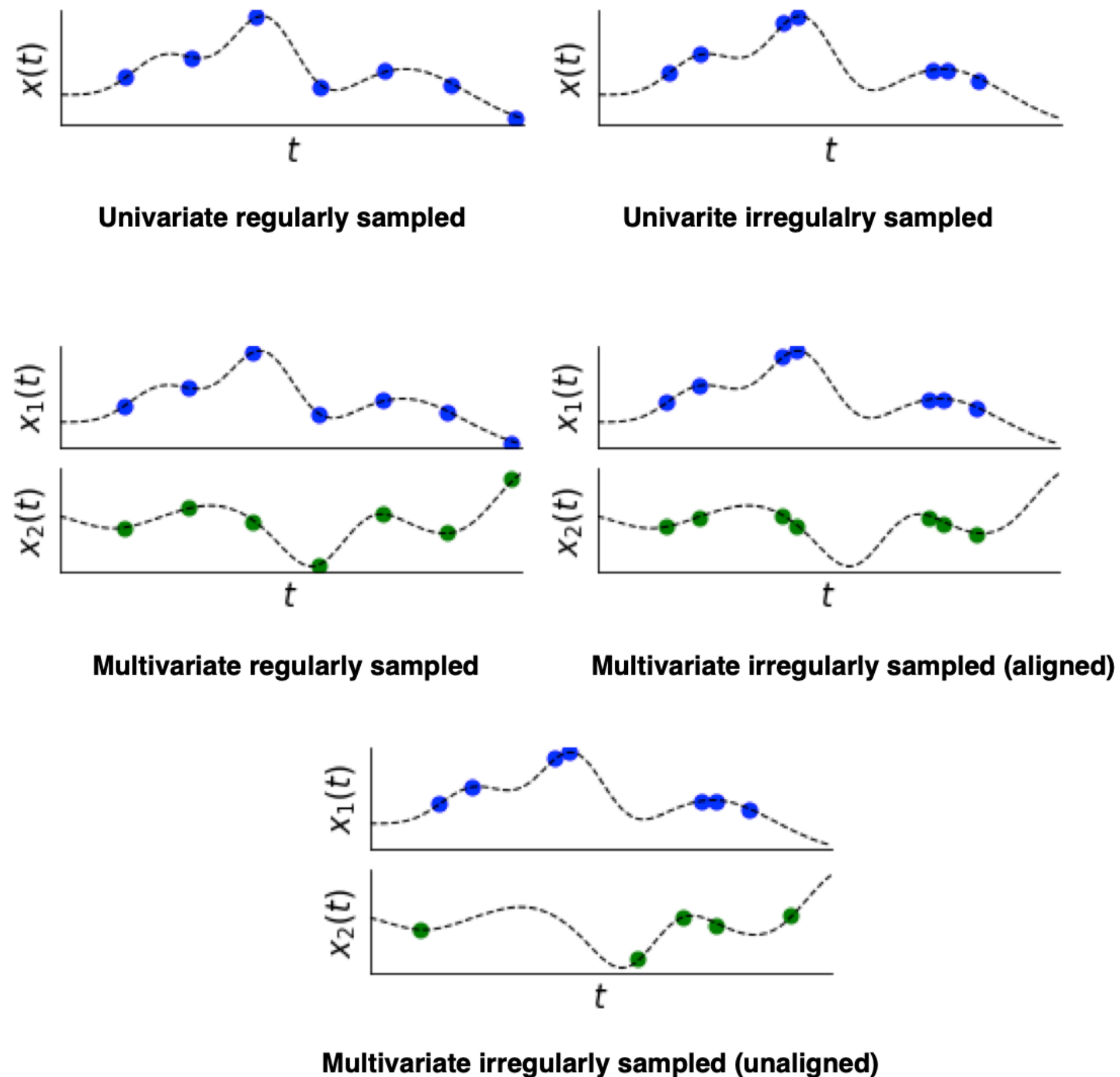
- ADC and associated analog amplifier may be power hungry
- Lossy compression may be power hungry (doing FFT, DCT etc.)
- In compressive sampling, the highest cost is random number generation
- To maximize benefit, may also need to duty cycle analog circuitry





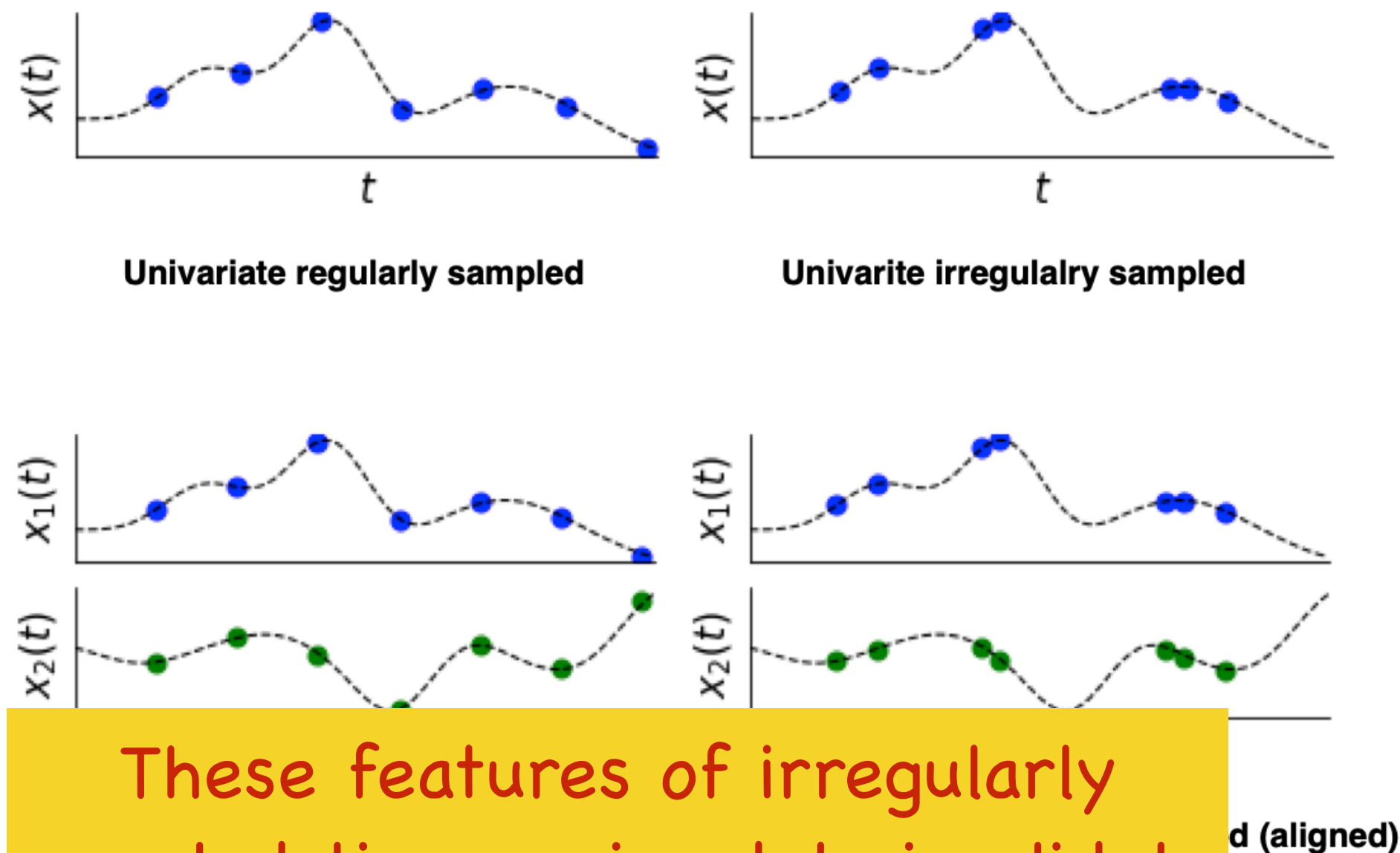
# Learning from Irregularly Sampled Time Series

# Irregularly Sampled Time Series Data



- Present fundamental challenges to many classical models from machine learning and statistics
- Consider supervised learning task where a model takes as input an irregularly sampled time series and must predict a scalar output
  - ▶ Training set  $\mathcal{D}$  with samples  $(s_i, y_i)$  where  $s_i$ 's are irregularly samples time series, and  $y_i$ 's are corresponding labels
- Problems:
  - ▶ **Variable gaps** between successive observation time points
  - ▶ **Variable number of observations:** the total # of observations across dimensions can vary across samples
  - ▶ **Lack of alignment:** Different dimensions of a single multivariate time series can be observed at a different collection of time points. The collection of observation times across dimensions can also differ between samples.

# Irregularly Sampled Time Series Data



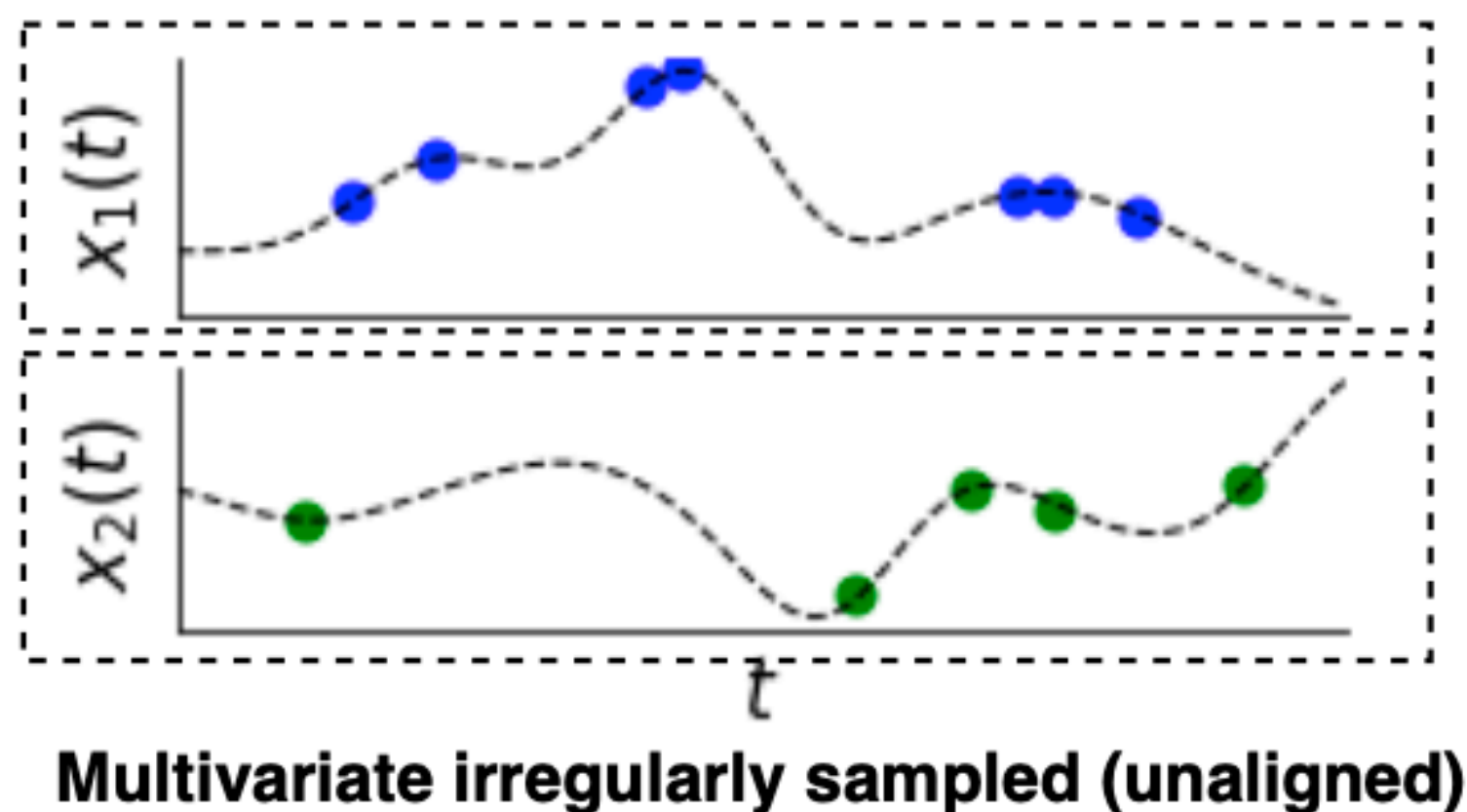
These features of irregularly sampled time series data invalidate the assumption of a coherent fixed-dimensional feature space, which underlies most basic supervised and un-supervised learning models.

- Present fundamental challenges to many classical models from machine learning and statistics
- Consider supervised learning task where a model takes as input an irregularly sampled time series and must predict a scalar output
  - ▶ Training set  $\mathcal{D}$  with samples  $(s_i, y_i)$  where  $s_i$ 's are irregularly sampled time series, and  $y_i$ 's are corresponding labels
- Problems:
  - ▶ **Variable gaps** between successive observation time points
  - ▶ **Variable number of observations**: the total # of observations across dimensions can vary across samples
  - ▶ **Lack of alignment**: Different dimensions of a single multivariate time series can be observed at a different collection of time points. The collection of observation times across dimensions can also differ between samples.

# Data Representations for Irregularly Sampled Time Series

- There are several possible data representations for multivariate irregularly sampled time series.
- These representations are equivalent, but expose different properties and suggest different approaches to modeling

## Series-based Representation



$$\begin{aligned}\mathbf{s} &= [\mathbf{s}_1, \mathbf{s}_2] \\ \mathbf{s}_1 &= (\mathbf{t}_1, \mathbf{x}_1) \\ \mathbf{t}_1 &= [t_{11}, \dots, t_{L_1 1}] \\ \mathbf{x}_1 &= [x_{11}, \dots, x_{L_1 1}] \\ \mathbf{s}_2 &= (\mathbf{t}_2, \mathbf{x}_2) \\ \mathbf{t}_2 &= [t_{12}, \dots, t_{L_2 2}] \\ \mathbf{x}_2 &= [x_{12}, \dots, x_{L_2 2}]\end{aligned}$$

A  $D$ -dimensional multivariate irregularly sampled time series  $\mathbf{s} = [\mathbf{s}_1, \dots, \mathbf{s}_D]$  is represented as a collection of univariate irregularly sampled time series, one per dimension.

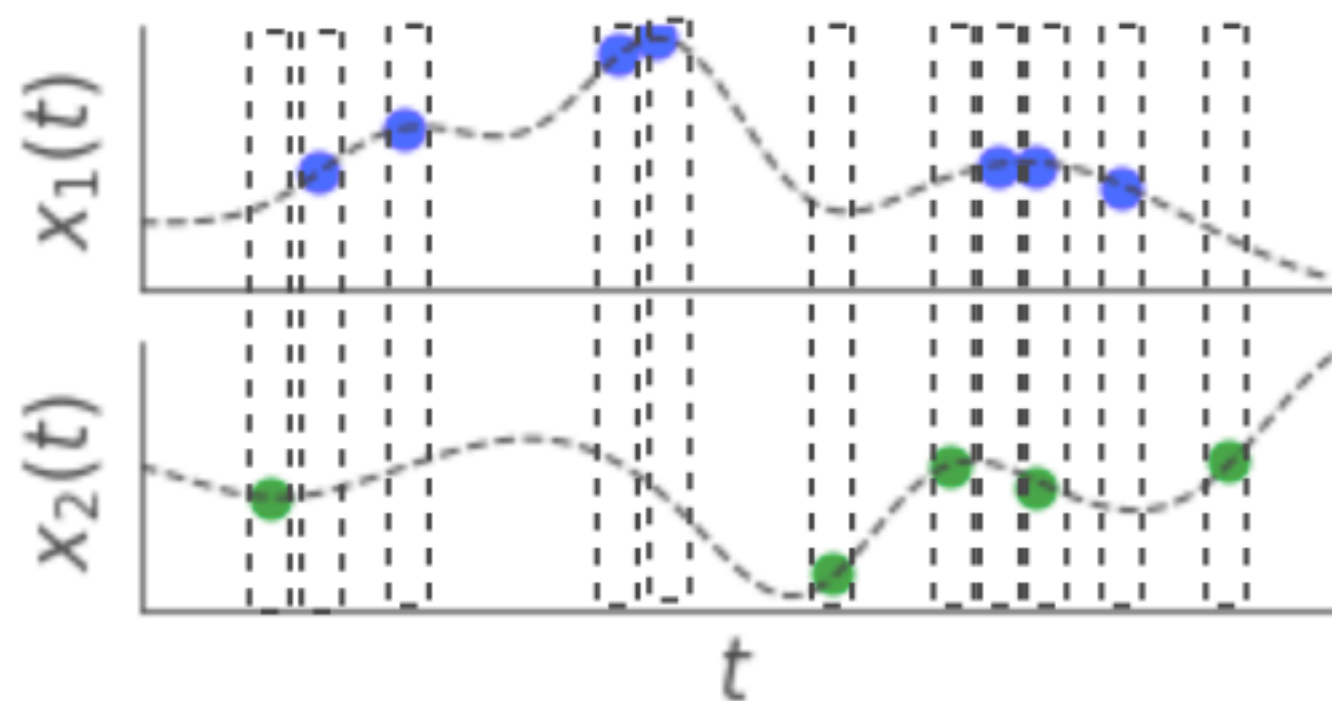
- $\mathbf{s}_d = (\mathbf{t}_d, \mathbf{x}_d)$  indicates the time series for dimension  $d$ .
- $\mathbf{t}_d$  indicates the collection of time points with observed values for dimension  $d$ .
- $\mathbf{x}_d$  indicates the corresponding collection of observed values.



# Data Representations for Irregularly Sampled Time Series

- There are several possible data representations for multivariate irregularly sampled time series.
- These representations are equivalent, but expose different properties and suggest different approaches to modeling

## Vector-based Representation



Multivariate irregularly sampled (unaligned)

$$\begin{aligned} \mathbf{s} &= (\mathbf{t}, \mathbf{x}, \mathbf{r}) \\ \mathbf{t} &= [t_1, \dots, t_L] \\ \mathbf{x} &= [\mathbf{x}_1, \dots, \mathbf{x}_L] \\ \mathbf{r} &= [\mathbf{r}_1, \dots, \mathbf{r}_L] \\ \mathbf{x}_i &= \begin{bmatrix} x_{i1} \\ x_{i2} \end{bmatrix} & x_{id} = \begin{cases} x_{id}, & \text{if } x_{id} \text{ is observed} \\ \text{NA}, & \text{otherwise} \end{cases} \\ \mathbf{r}_i &= \begin{bmatrix} r_{i1} \\ r_{i2} \end{bmatrix} & r_{id} = \begin{cases} 1, & \text{if } x_{id} \text{ is observed} \\ 0, & \text{otherwise} \end{cases} \end{aligned}$$

In this representation, there is a single collection of time points  $\mathbf{t}$ . At each time point  $t_i$ , there is a  $D$ -dimensional vector-valued observation  $\mathbf{x}_i$ .

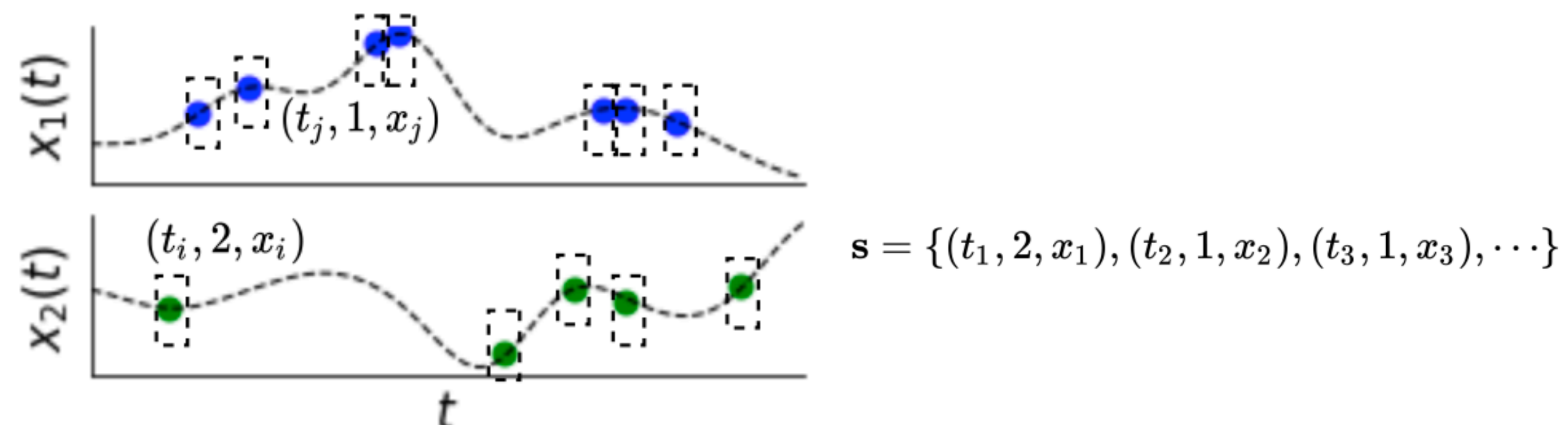
In the general case, not all dimensions of  $\mathbf{x}_i$  are observed, leading to the need to explicitly represent which dimensions are observed and which are missing.

A  $D$ -dimensional binary response indicator vector  $\mathbf{r}_i$  at each time point  $t_i$  indicates which dimensions are observed and which are missing.

# Data Representations for Irregularly Sampled Time Series

- There are several possible data representations for multivariate irregularly sampled time series.
- These representations are equivalent, but expose different properties and suggest different approaches to modeling

## Set-based Representation



**Multivariate irregularly sampled (unaligned)**

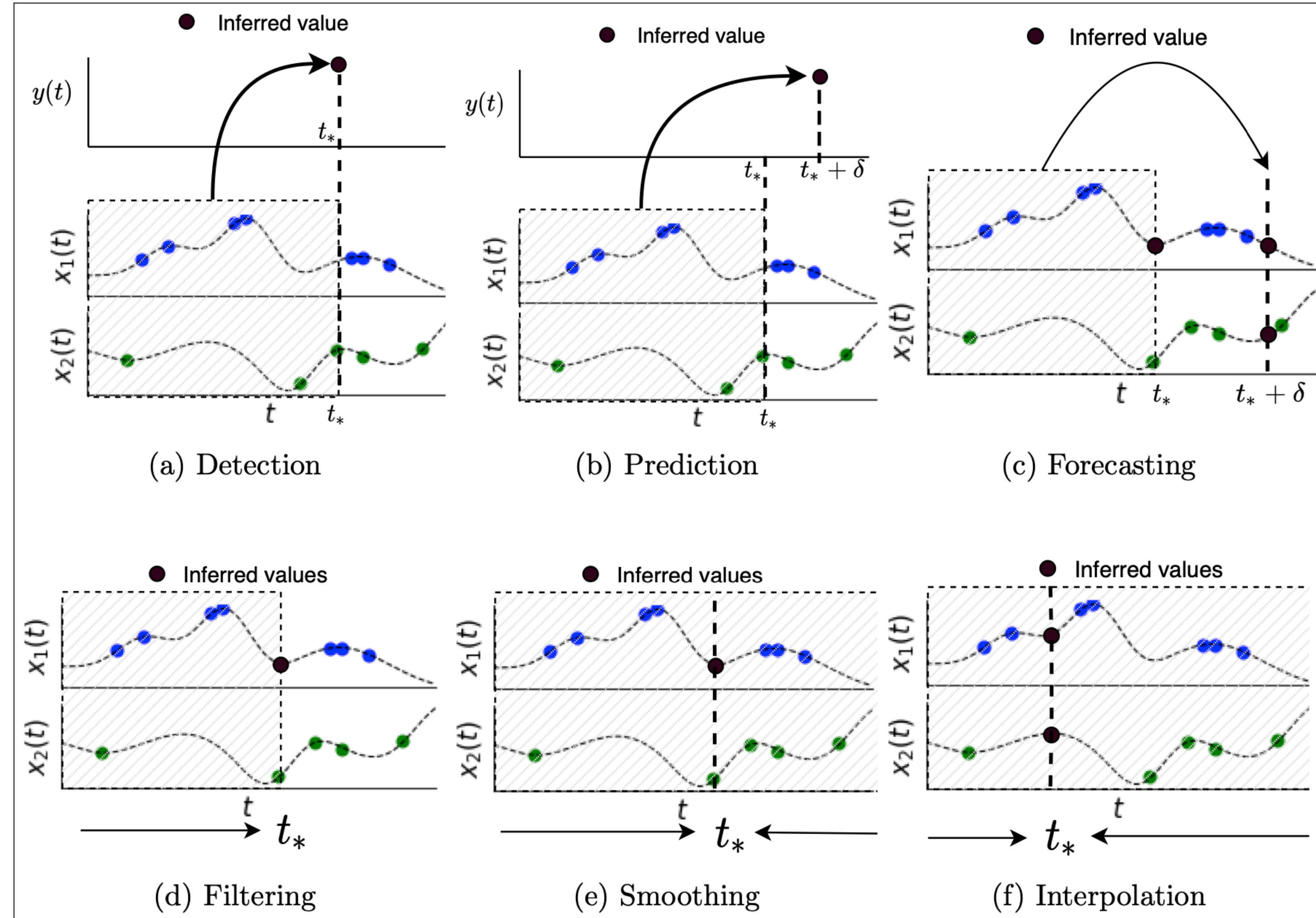
In this representation, a  $D$ -dimensional multivariate irregularly sampled time series is represented as a set of *(time, dimension, value)* tuples, one for each observation.



# Inference Tasks

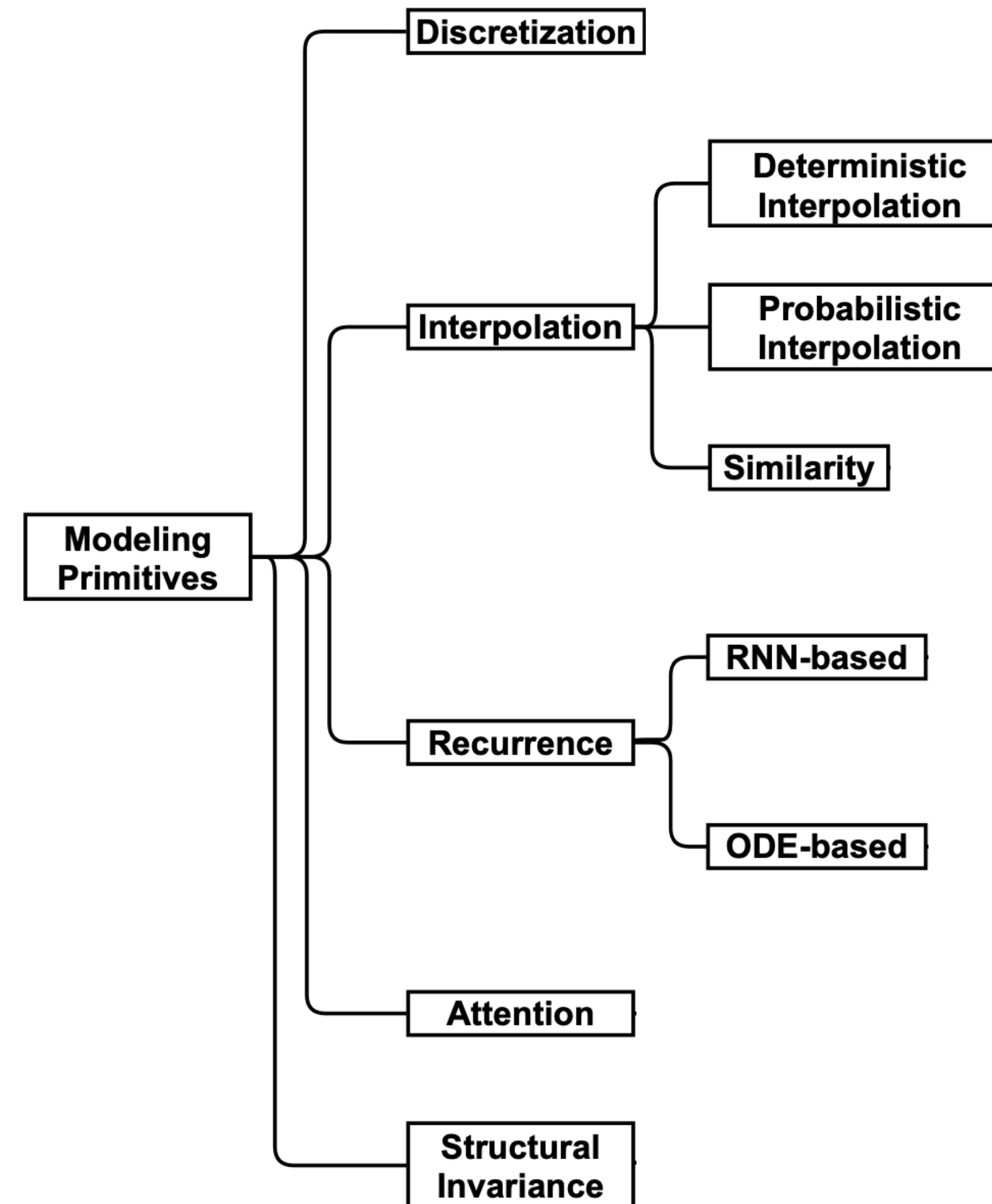
Note that in machine learning, the inference for any quantity that is not known is often referred to as a “prediction”, but here we use the term “prediction” to refer to a task where the inference is for the value of the output variable at a time that is in the relative future of the time point  $t_*$  at which the inference is made.

- **Detection:** Inferring prediction target values  $y[t_*]$  at time  $t_*$  conditioning on the observations  $\mathbf{s}[: t_*]$  available up to and including time  $t_*$ .
- **Prediction:** Inferring prediction target values  $y[t_* + \delta]$  at time  $t_* + \delta$  (for  $\delta > 0$ ) conditioning on the observations  $\mathbf{s}[: t_*]$  available up to and including time  $t_*$ .
- **Forecasting:** Inferring  $\mathbf{x}[t_* + \delta]$  (for  $\delta > 0$ ) by conditioning on the observations  $\mathbf{s}[: t_*]$  up to and including time  $t_*$ .
- **Filtering:** Inferring missing variables  $\mathbf{x}^m[t_*]$  at time  $t_*$  by conditioning on the observations  $\mathbf{s}[: t_*]$  up to and including time  $t_*$ .
- **Smoothing:** Inferring the values of  $\mathbf{x}^m[t_*]$  at time  $t_*$  using the observed data in  $\mathbf{s}$ .
- **Interpolation:** Inferring the values of  $\mathbf{x}[t_*]$  at time  $t_*$  using the observed data in  $\mathbf{s}$ .



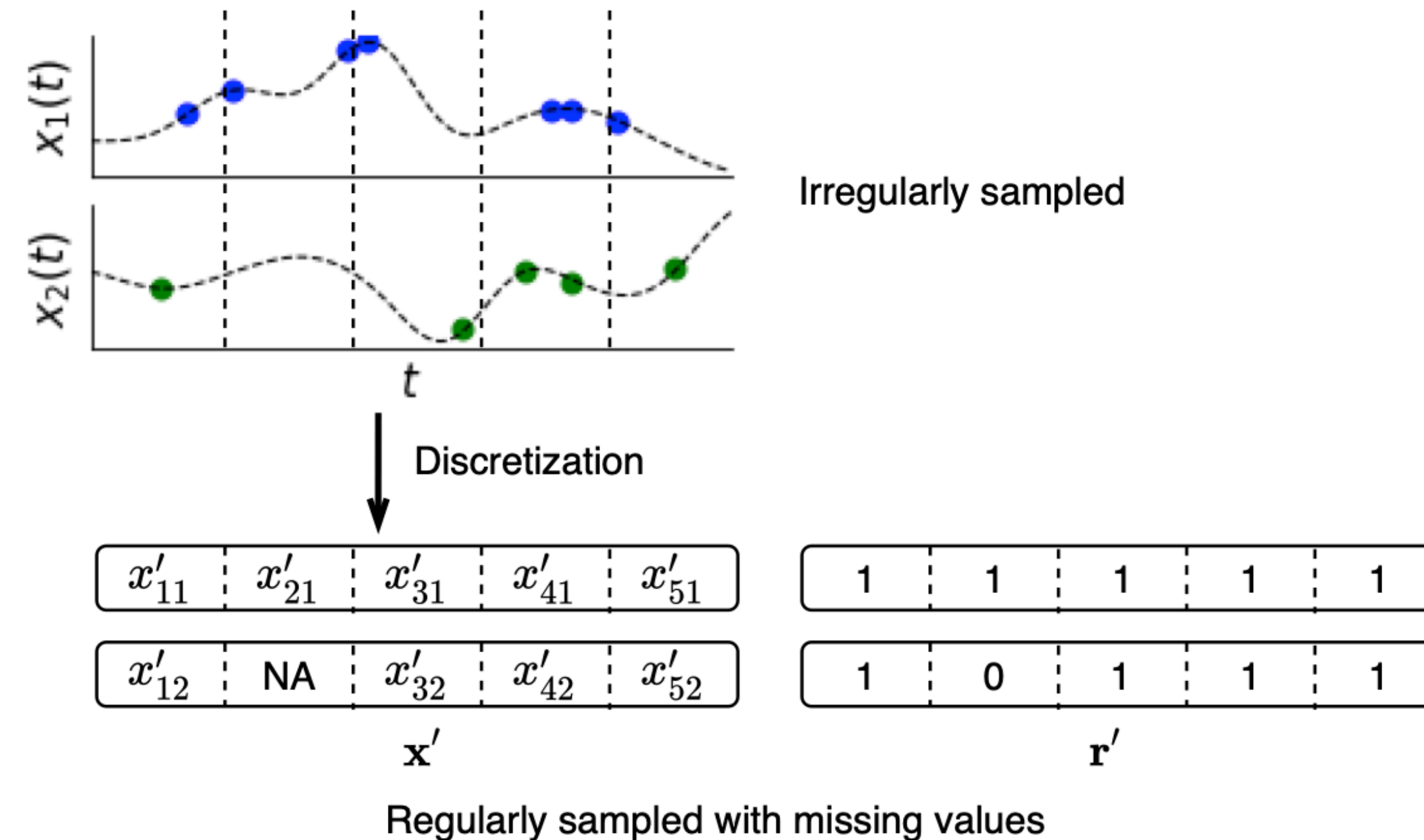
# Modeling Primitives for Irregularly Sampled Time Series

---



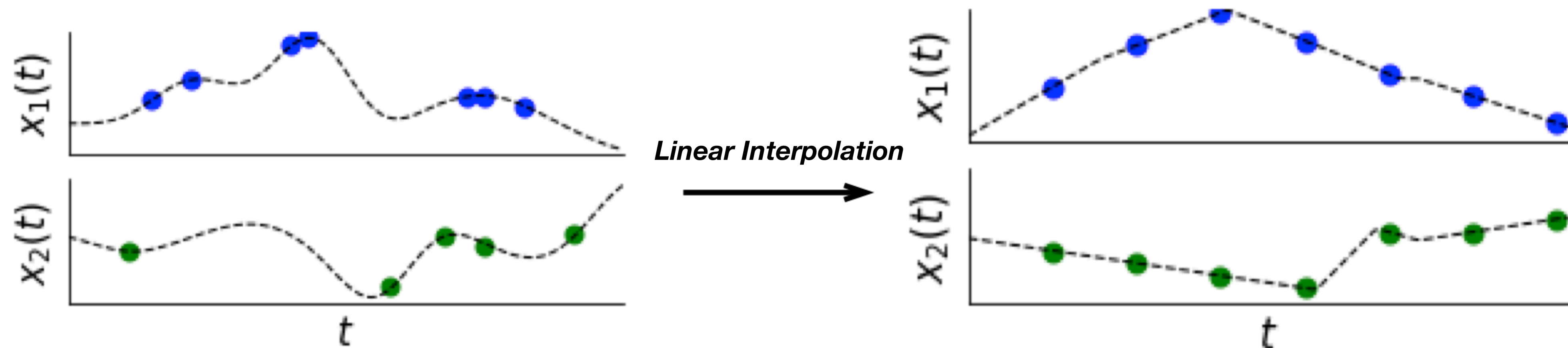


# Discretization



- Reduces to a regularly sampled multivariate time series with missing values
  - vector-based representation with missing data indicator
- Approach: divide the time axis into equal sized non-overlapping intervals and define a value within each time interval based on the observed values falling within that interval
  - e.g. average or median

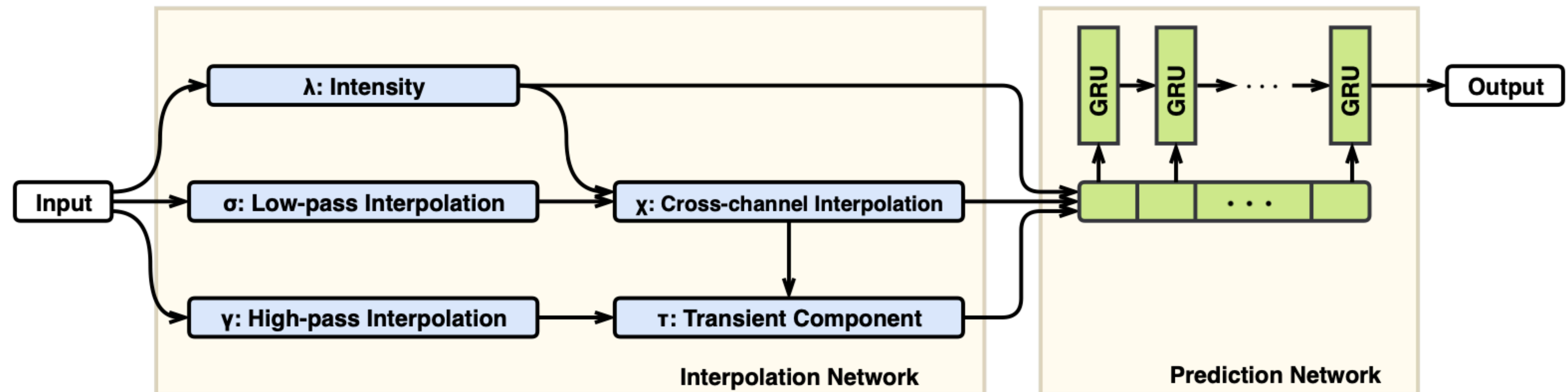
# Interpolation



- Approach:
  1. Define a set of  $K$  reference time points  $\tau = [\tau_1, \dots, \tau_K]$
  2. Use a basic kernel smoother to produce interpolated values at the reference time points
    - kernel typically puts higher weights (learnable parameters) to points that are closer to the reference points
    - for multivariate case, can account for both correlation in time and correlation across different dimensions
    - deterministic: e.g., squared exponential kernel, stochastic: e.g., Gaussian process regression

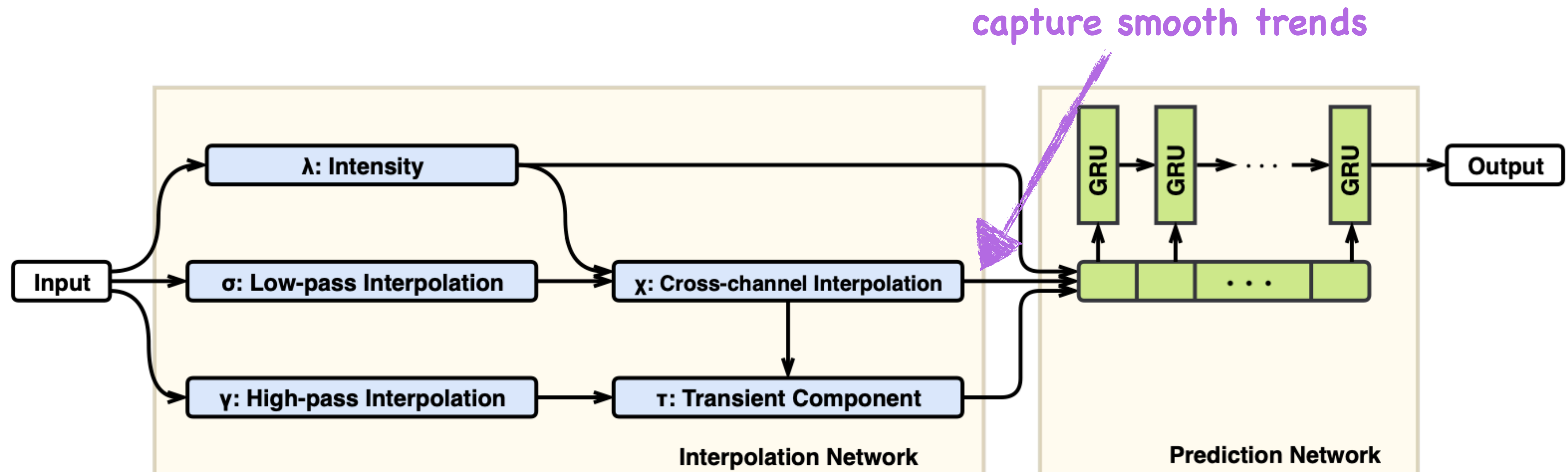
# Example: Interpolation-Prediction Networks (IP-Net)

- Applies multiple semi-parametric interpolation schemes to obtain a regularly-sampled time series representation with samples at a set of reference time points
  - The parameters of the interpolation network are trained with the classifier in an end-to-end setup
- Prediction Network can be any standard supervised neural network architecture such as fully-connected feedforward, convolutional, or recurrent network.



# Example: Interpolation-Prediction Networks (IP-Net)

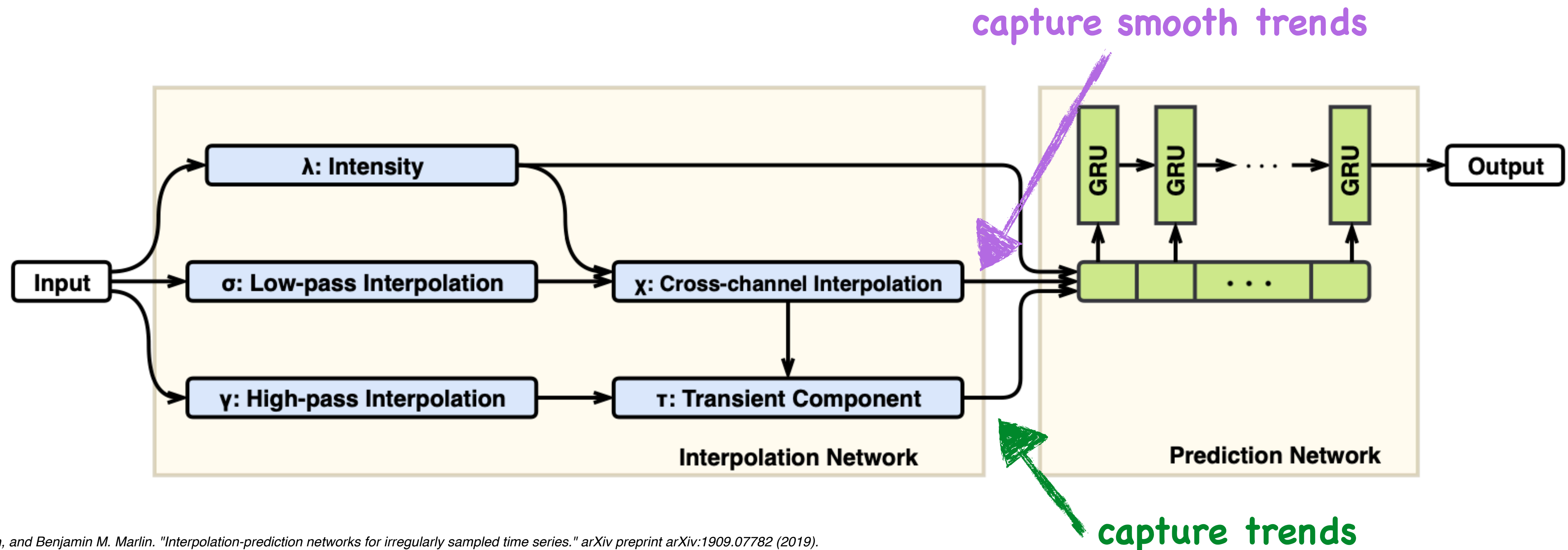
- Applies multiple semi-parametric interpolation schemes to obtain a regularly-sampled time series representation with samples at a set of reference time points
  - The parameters of the interpolation network are trained with the classifier in an end-to-end setup
- Prediction Network can be any standard supervised neural network architecture such as fully-connected feedforward, convolutional, or recurrent network.





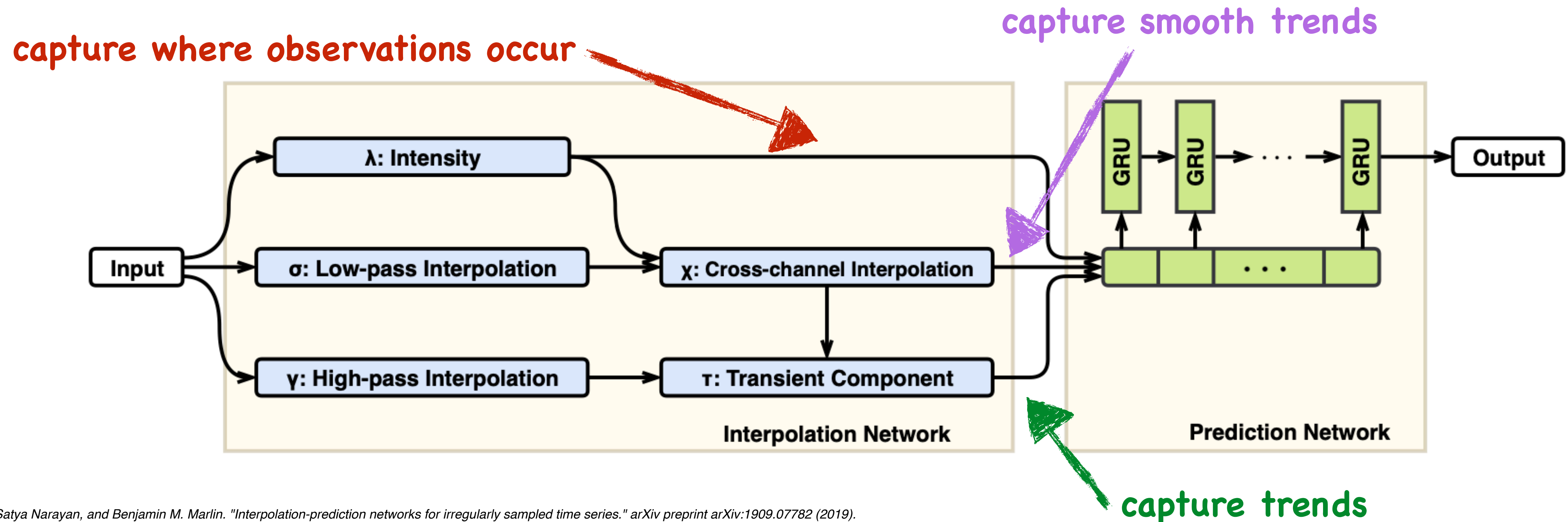
# Example: Interpolation-Prediction Networks (IP-Net)

- Applies multiple semi-parametric interpolation schemes to obtain a regularly-sampled time series representation with samples at a set of reference time points
  - The parameters of the interpolation network are trained with the classifier in an end-to-end setup
- Prediction Network can be any standard supervised neural network architecture such as fully-connected feedforward, convolutional, or recurrent network.



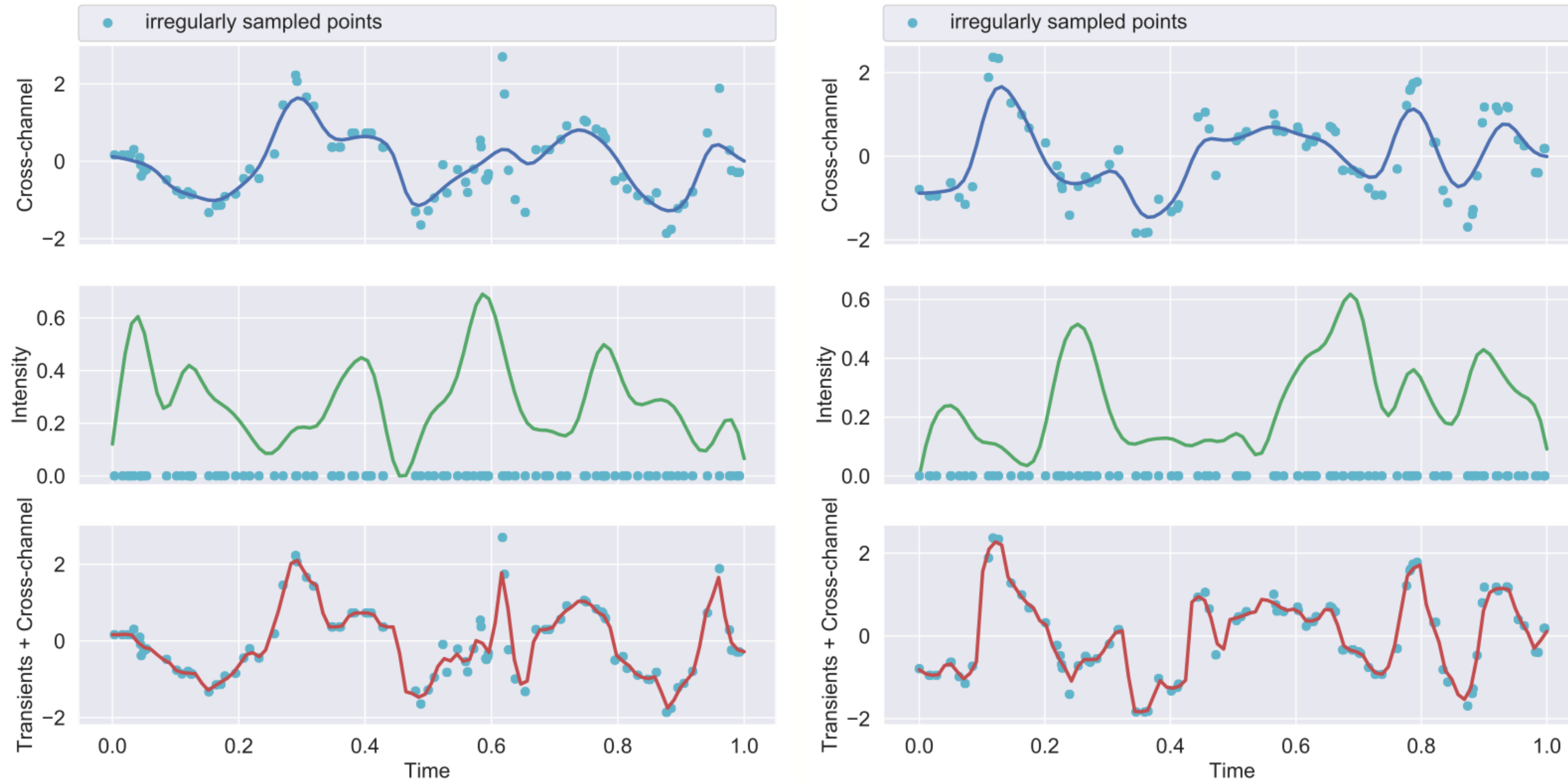
# Example: Interpolation-Prediction Networks (IP-Net)

- Applies multiple semi-parametric interpolation schemes to obtain a regularly-sampled time series representation with samples at a set of reference time points
  - The parameters of the interpolation network are trained with the classifier in an end-to-end setup
- Prediction Network can be any standard supervised neural network architecture such as fully-connected feedforward, convolutional, or recurrent network.

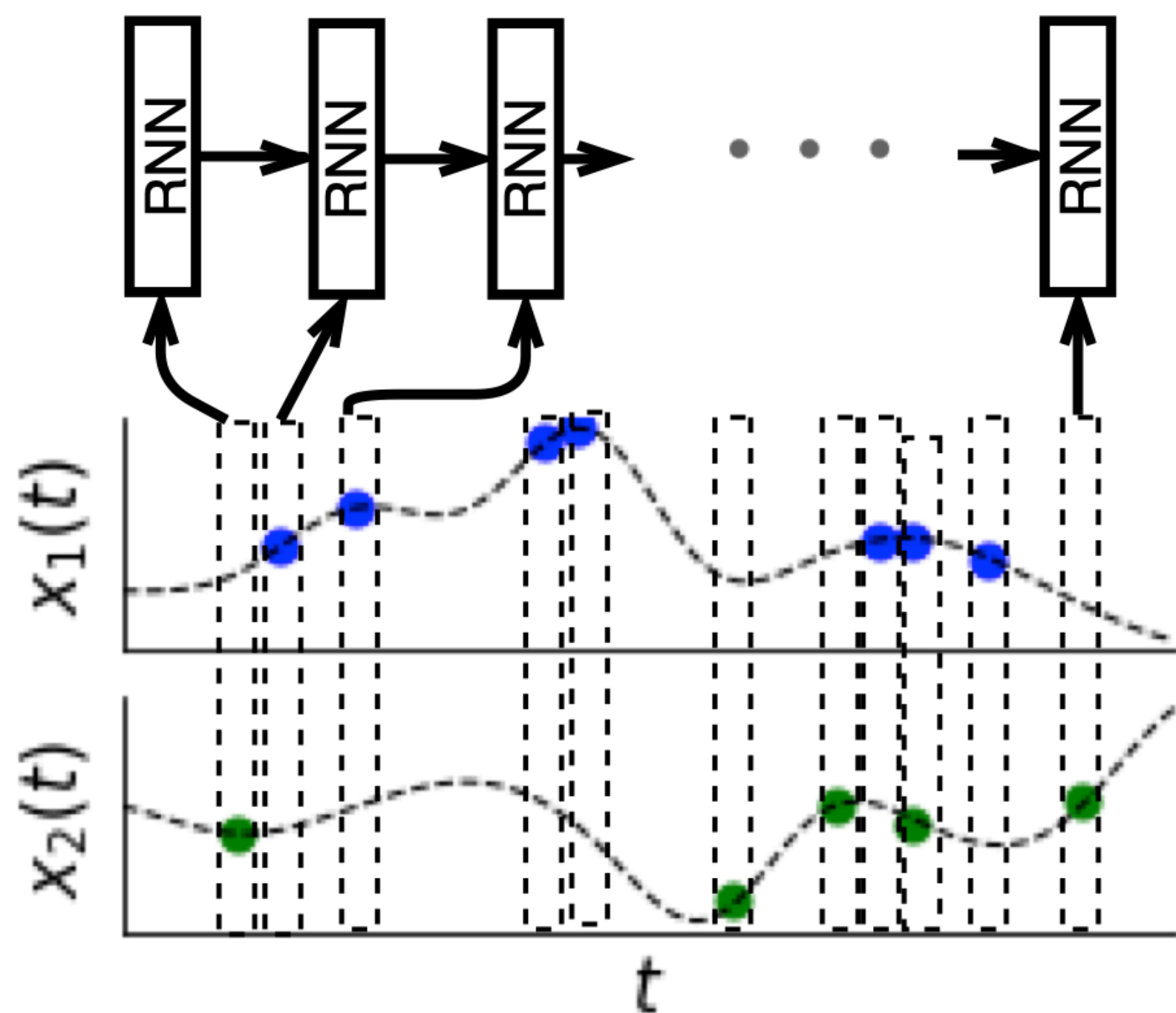




# Output of Interpolation Network



# Recurrence



**Multivariate irregularly sampled (unaligned)**

- Use a RNN cell with the ability to explicitly represent time to integrate the input at each time point with the latent state from the previous time point
  - ▶ e.g. append the time points or inter-observation intervals to the vector-valued observations  $[\mathbf{x}_{in}, t_{in}]$  or  $[\mathbf{x}_{in}, t_{in} - t_{i-1n}]$

$$\mathbf{h}_i = f_{\theta}(\mathbf{h}_{i-1}, \mathbf{z}_{in})$$

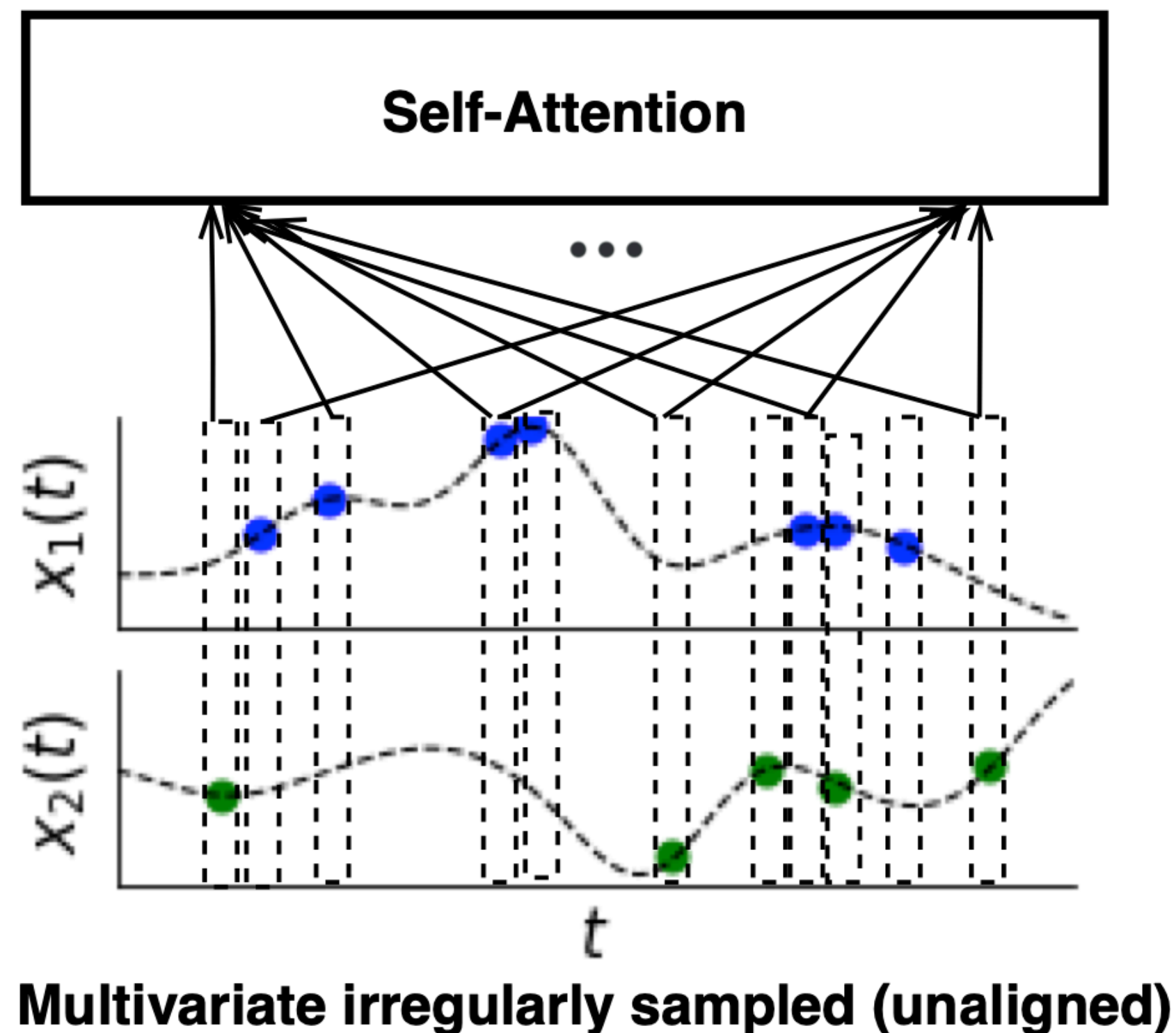
$$\hat{\mathbf{y}}_{in} = q_{\phi}(\mathbf{h}_i)$$

- Recent work on ordinary differential equation (ODE) models in ML provides an alternative recurrence-based solution
  - ▶ In these ODE-RNN models, ODEs are used to evolve the hidden state between continuous time observations.
  - ▶ Better properties than traditional RNNs in terms of their ability to accommodate irregularly sampled data.

$$\mathbf{h}'_i = \text{ODESolve}(g_{\gamma}, \mathbf{h}_{i-1}, (t_{i-1n}, t_{in}))$$

$$\mathbf{h}_i = f_{\theta}(\mathbf{h}'_i, \mathbf{x}_{in})$$

# Attention



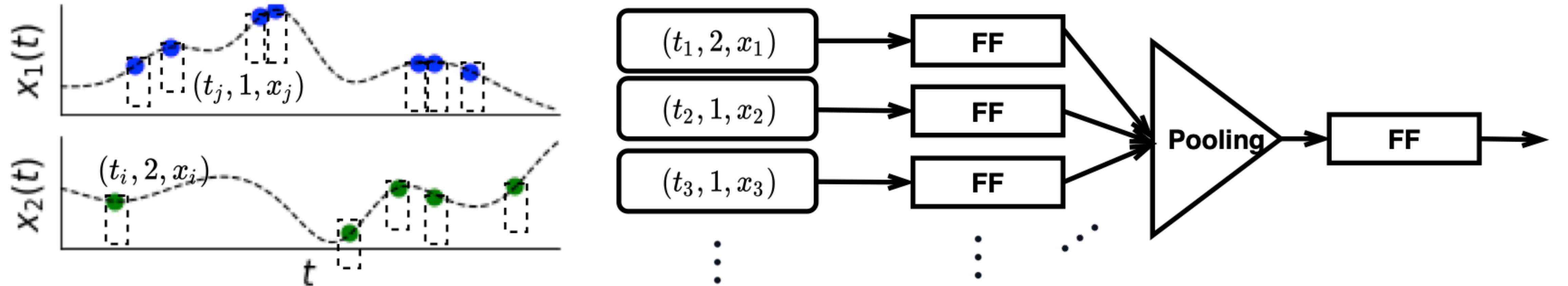
- Self-Attention module learns which regions of an input time series to attend to when computing outputs at different points in time by leveraging positional or time encodings.

$$\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{C}}\right)\mathbf{V}$$

- There is no recurrent structure, which allows parallel processing go the entire time series instead of sequentially as in RNN
- Time values can be converted into a vector representation using positional encoding and concatenated with the observation value (as in RNN)
- Missing values in vector-valued observations are also problematic for attention-based modules, which (like standard RNNs) expect fully observed vectors as input
  - imputation solutions can be used



# Structural Invariance



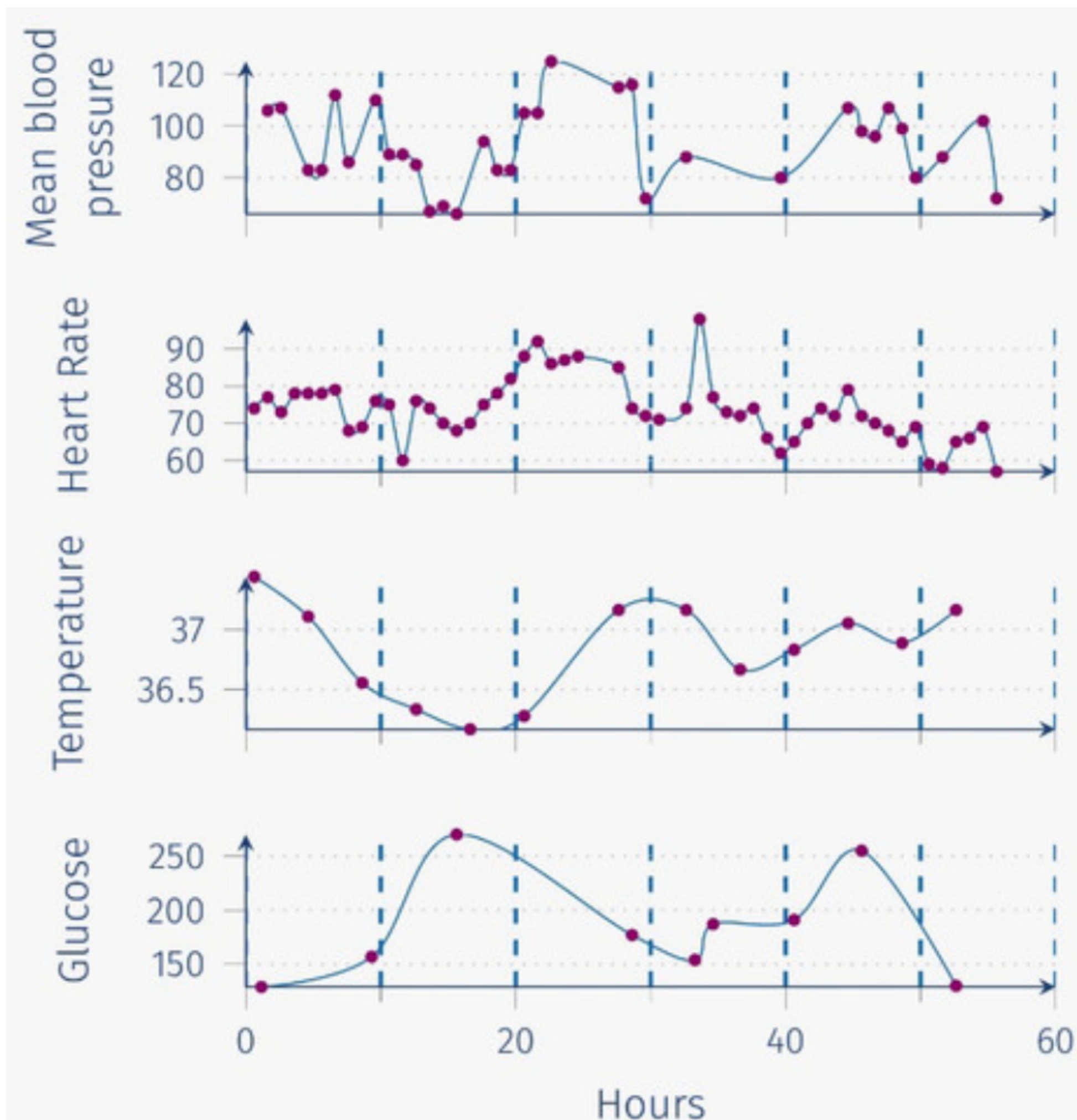
**Multivariate irregularly sampled (unaligned)**

- Inspired by the set-based view of a multivariate irregularly sampled time series  $\mathbf{s}_n = (t_{in}, d_{in}, \mathbf{x}_{in}) \mid 1 \leq i \leq L_n$
- Set-based neural network approach processes individual (time, value, dimension) tuples via an encoding function and then pools over the output of all such tuples
  - such approaches (i) produce an encoding of an input irregularly sampled time series by applying an initial encoder  $f_\theta$  to individual *time-dimension-value* tuples, (ii) then perform a pooling operation (e.g. max, mean and sum) over all initial encodings in a way that is completely invariant to the temporal structure of the data, and (iii) map the output of pooling through one additional set of encoding layers  $g_\phi$  to produce the final representation

$$\mathbf{h} = g_\phi(\text{pool}(f_\theta(t_{in}, d_{in}, \mathbf{x}_{in}) \mid 1 \leq i \leq L_n))$$

# Example Approach: Set Functions for Time Series (SeFT)

## Medical Time Series



### Challenge

Imputation requires solving the harder problem of learning the time series dynamics, and also sacrifices interpretability of the inference.

### Problem Statement

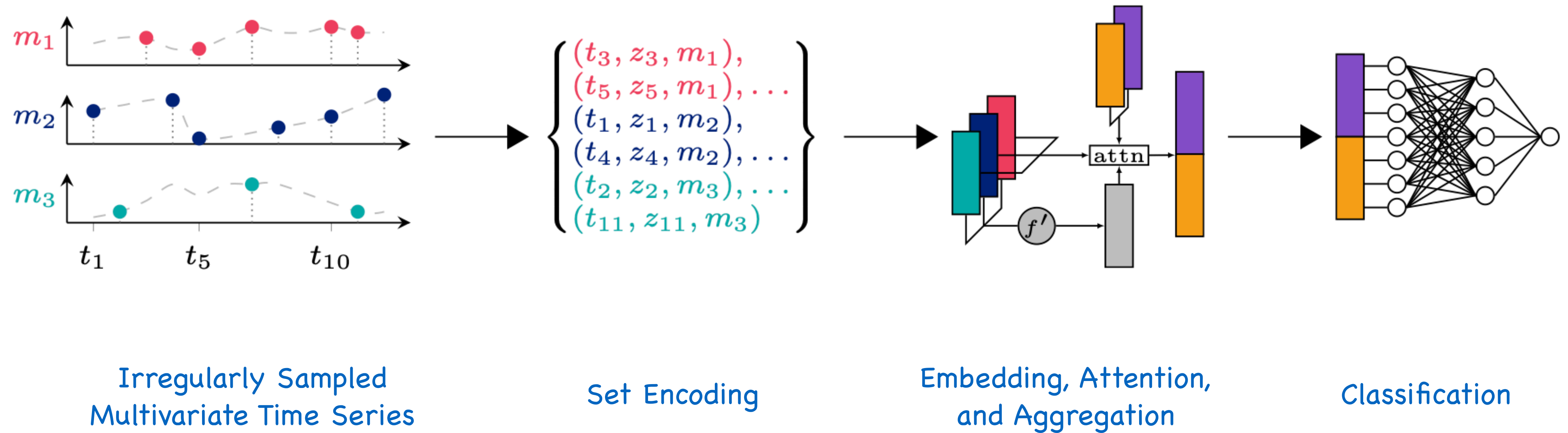
Can we learn classification models on irregularly sampled time series without prior imputation?

### Set Functions for Time Series

→ Time series classification as set classification

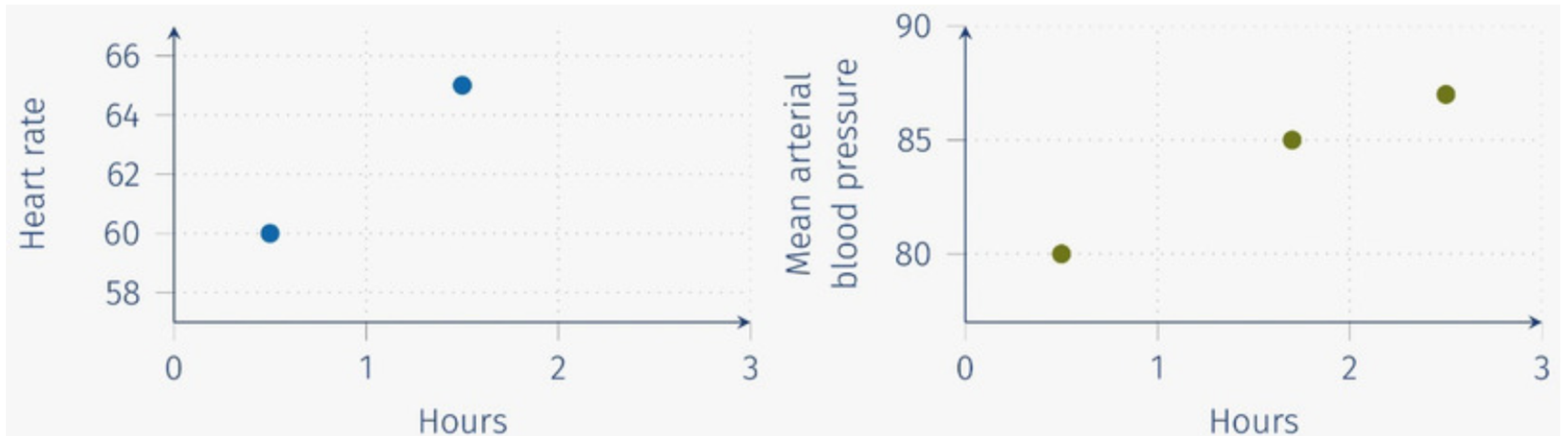
(based on recent advances in differentiable set function learning)

# SeFT Architecture Overview





# Key Idea in SeFT: Time Series as Set of Observations



Each observation  $s_j$  is represented as a tuple  $(t_j, z_j, m_j)$

$$\mathcal{S} = \{(0.5, 60, 1), (1.5, 65, 1), (0.5, 80, 2), (1.7, 85, 2), (3, 87, 2)\}$$

# Background: *Deep Sets* Framework (Zaheer et. al., NeurIPS 2017)

---

- Neural network architectures that operate over sets
  - ▶ encountered in many applications: compute statistics over sets (e.g. sum of digits in a set of images), classify sets (e.g. LIDAR or RADAR point cloud classification)
- Key requirement:
  - ▶ **permutation invariance** if the task is  $f : 2^{\mathcal{X}} \rightarrow \mathcal{Y}$  (i.e. regression or classification)
  - ▶ **permutation equivariance** of the task is  $f : \mathcal{X}^M \rightarrow \mathcal{Y}^M$  (i.e. transduction)
- Key Results

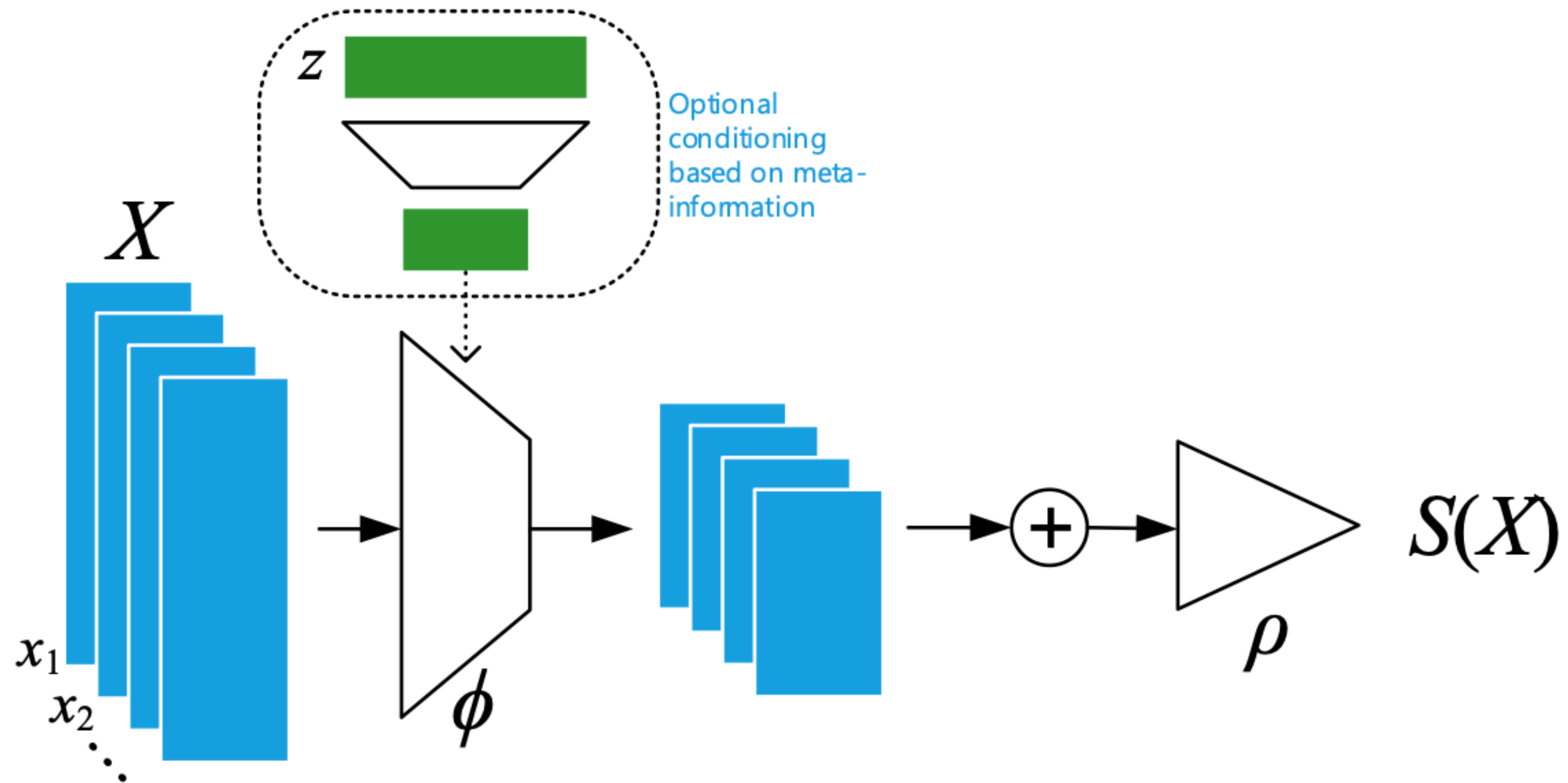
**Theorem 2** *A function  $f(X)$  operating on a set  $X$  having elements from a countable universe, is a valid set function, i.e., **invariant** to the permutation of instances in  $X$ , iff it can be decomposed in the form  $\rho \left( \sum_{x \in X} \phi(x) \right)$ , for suitable transformations  $\phi$  and  $\rho$ .*

$\mathbf{f}_{\Theta}(\mathbf{x}) = \sigma(\Theta \mathbf{x})$  is restricted to standard neural network layer

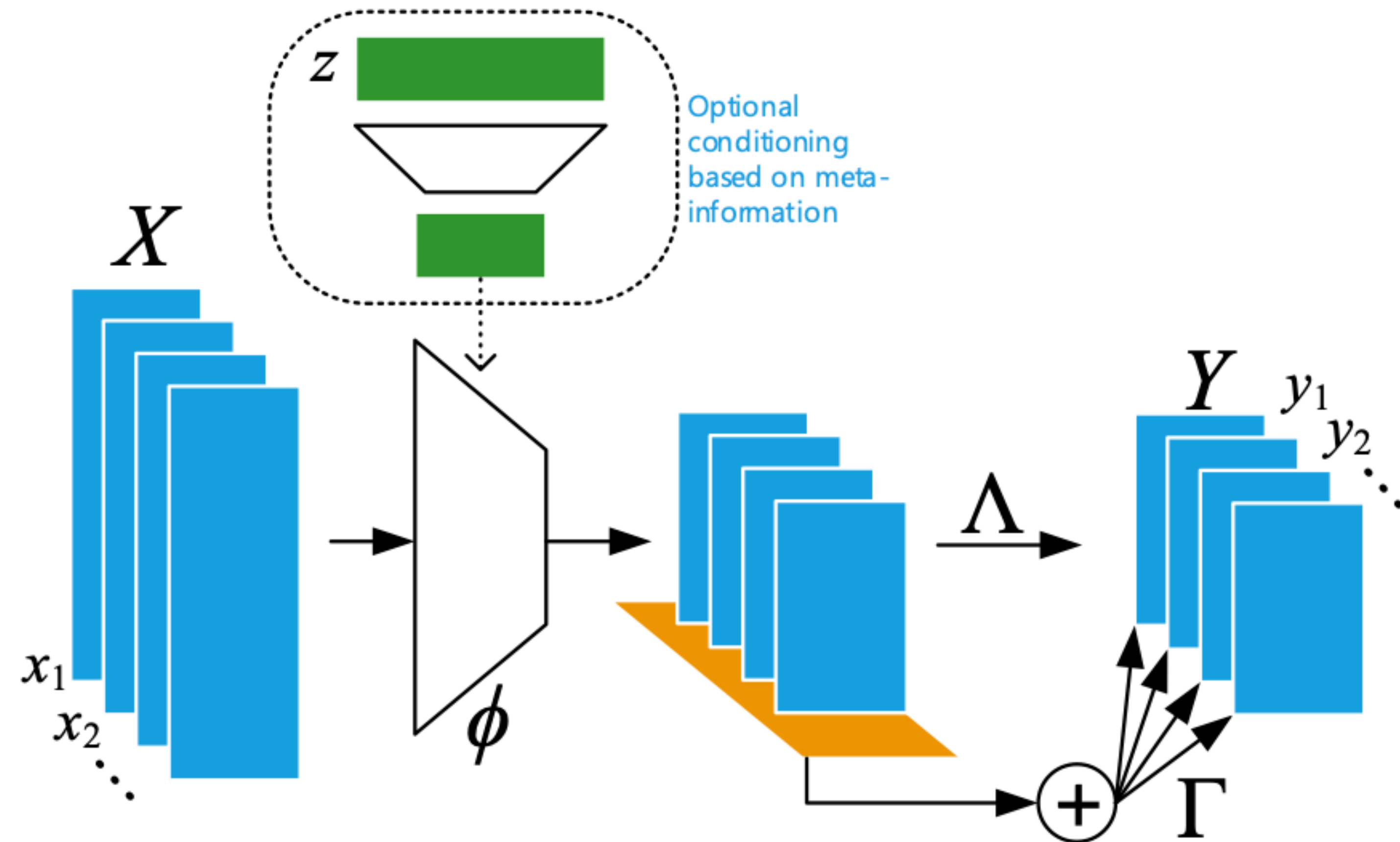
**Lemma 3** *The function  $\mathbf{f}_{\Theta} : \mathbb{R}^M \rightarrow \mathbb{R}^M$  defined above is permutation **equivariant** iff all the off-diagonal elements of  $\Theta$  are tied together and all the diagonal elements are equal as well. That is,*

$$\Theta = \lambda \mathbf{I} + \gamma (\mathbf{1}\mathbf{1}^T) \quad \lambda, \gamma \in \mathbb{R} \quad \mathbf{1} = [1, \dots, 1]^T \in \mathbb{R}^M \quad \mathbf{I} \in \mathbb{R}^{M \times M} \text{ is the identity matrix}$$

# Background: Architecture of DeepSets *Invariant* Model

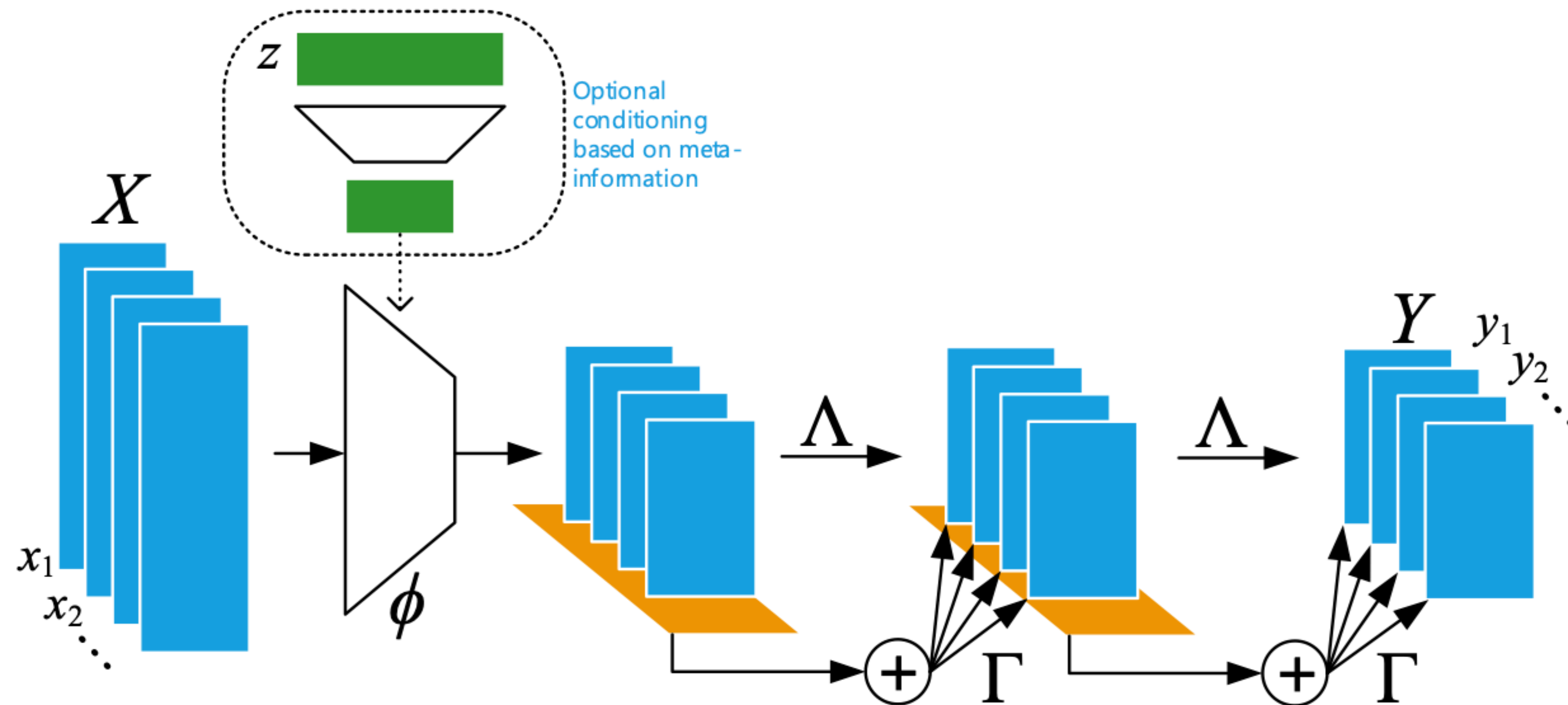


# Background: Architecture of DeepSets *Equivariant* Model





# Background: Architecture of DeepSets *Equivariant* Model



**Using multiple permutation equivariant layers.**  
(Since permutation equivariance compose we can stack multiple such layers.)

# Applying Deep Sets to SeFT

---

- ▶ Sum-decompose the set function to achieve permutation invariance

$$f(\mathcal{S}) = g \left( \frac{1}{|\mathcal{S}|} \sum_{s_j \in \mathcal{S}} h(s_j) \right)$$

where  $h: \Omega \rightarrow \mathbb{R}^d$  and  $g: \mathbb{R}^d \rightarrow \mathbb{R}^c$  are neural networks

## Problem

Influence of an element shrinks as  $|\mathcal{S}|$  grows!



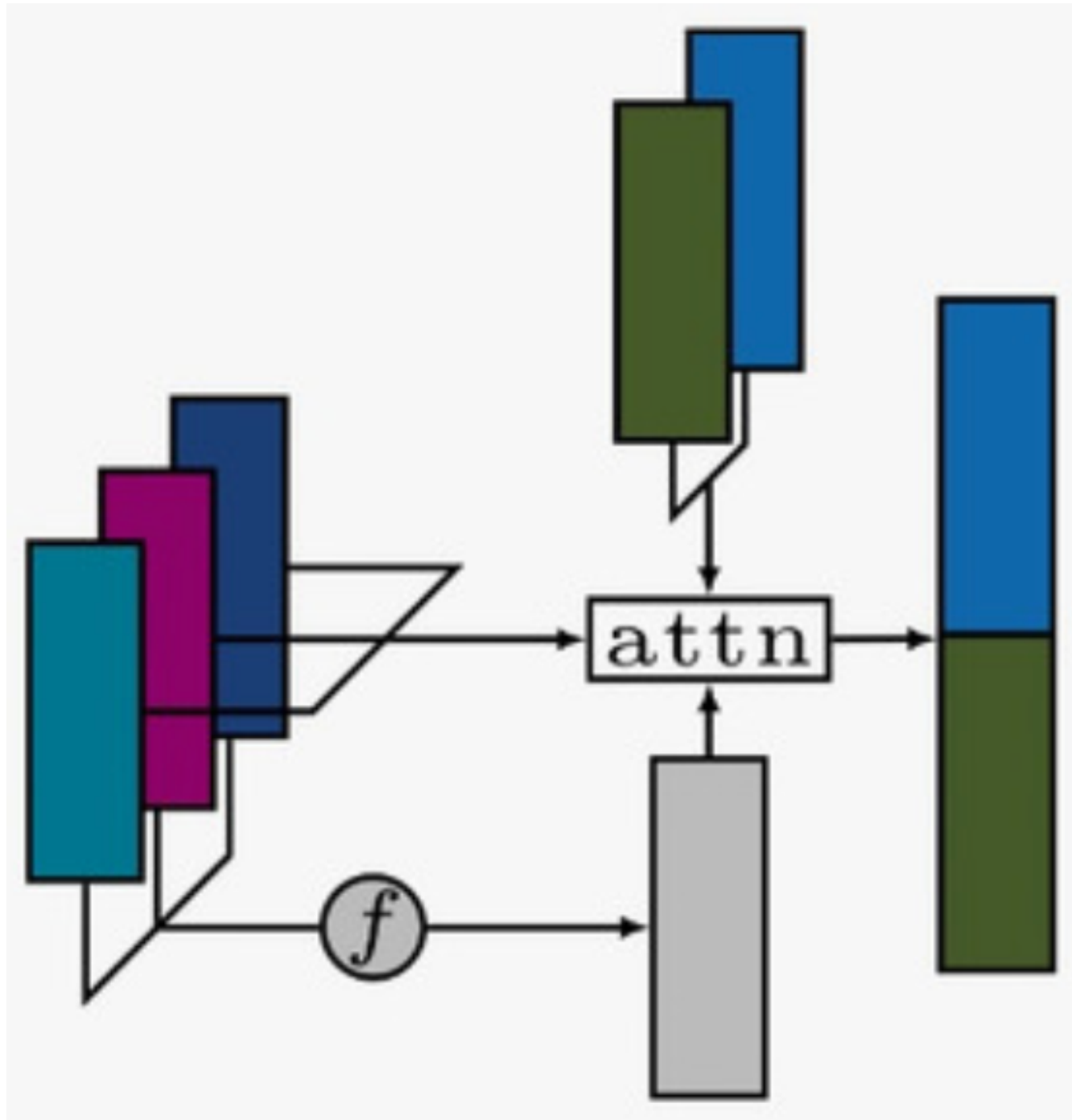
# Time Encoding in SeFT

---

- Employs a variant of *positional encoding* seen earlier with transformer
- The time encoding converts the 1-dimensional time axis into a multi-dimensional input by passing the time  $t$  of each observation through multiple trigonometric functions of varying frequencies
- Given a dimensionality  $\tau \in \mathbb{N}^+$  of the time encoding, SeFT encodes the position as  $x \in \mathbb{R}^\tau$ , where

$$\begin{aligned} x_{2k}(t) &:= \sin \left( \frac{t}{t^{2k/\tau}} \right) \\ x_{2k+1}(t) &:= \cos \left( \frac{t}{t^{2k/\tau}} \right) \end{aligned} \quad k \in \{0, \dots, \tau/2\}$$

# Attention-based Aggregation



$$\text{Keys: } K_{j,i} = [f(\mathcal{S}), s_j]^T W_i$$

$$\text{Queries: } Q \in \mathbb{R}^{m \times d}$$

$$\text{Preattentions: } e_{j,i} = \frac{K_{j,i} \cdot Q_i}{\sqrt{d}}$$

$$\text{Attentions: } a_{j,i} = \frac{\exp(e_{j,i})}{\sum_j \exp(e_{j,i})}$$

$$\text{Values: } V_i = \sum_j a_{j,i} h_{\theta}(s_j)$$

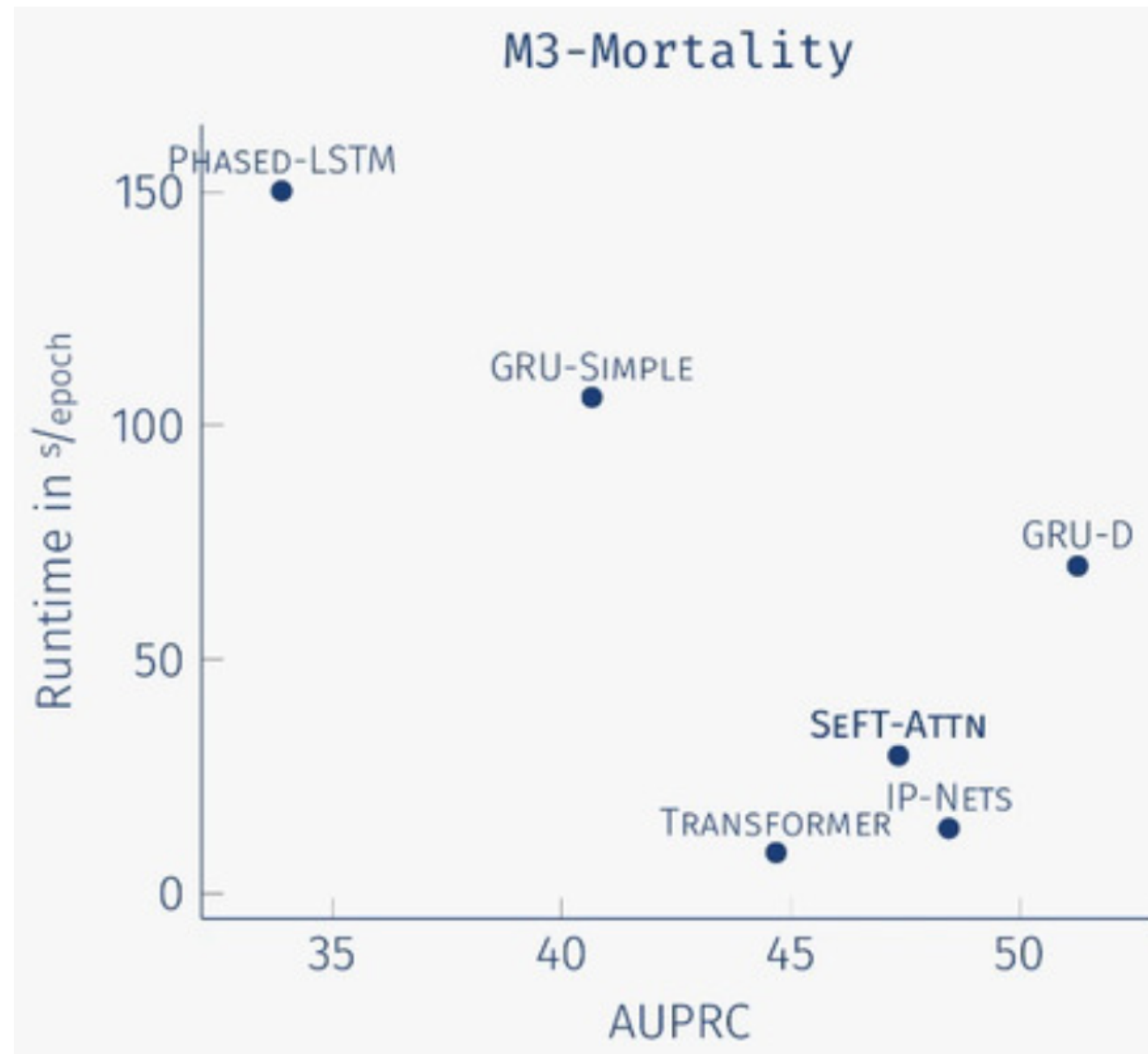
$$\mathcal{L}(\theta, \psi) = \mathbb{E}_{(\mathcal{S}, y) \in \mathcal{D}} \left[ \ell \left( y; g_{\psi} \left( \sum_{s_j \in \mathcal{S}} a(\mathcal{S}, s_j) h_{\theta}(s_j) \right) \right) \right]$$

# SeFT Evaluation

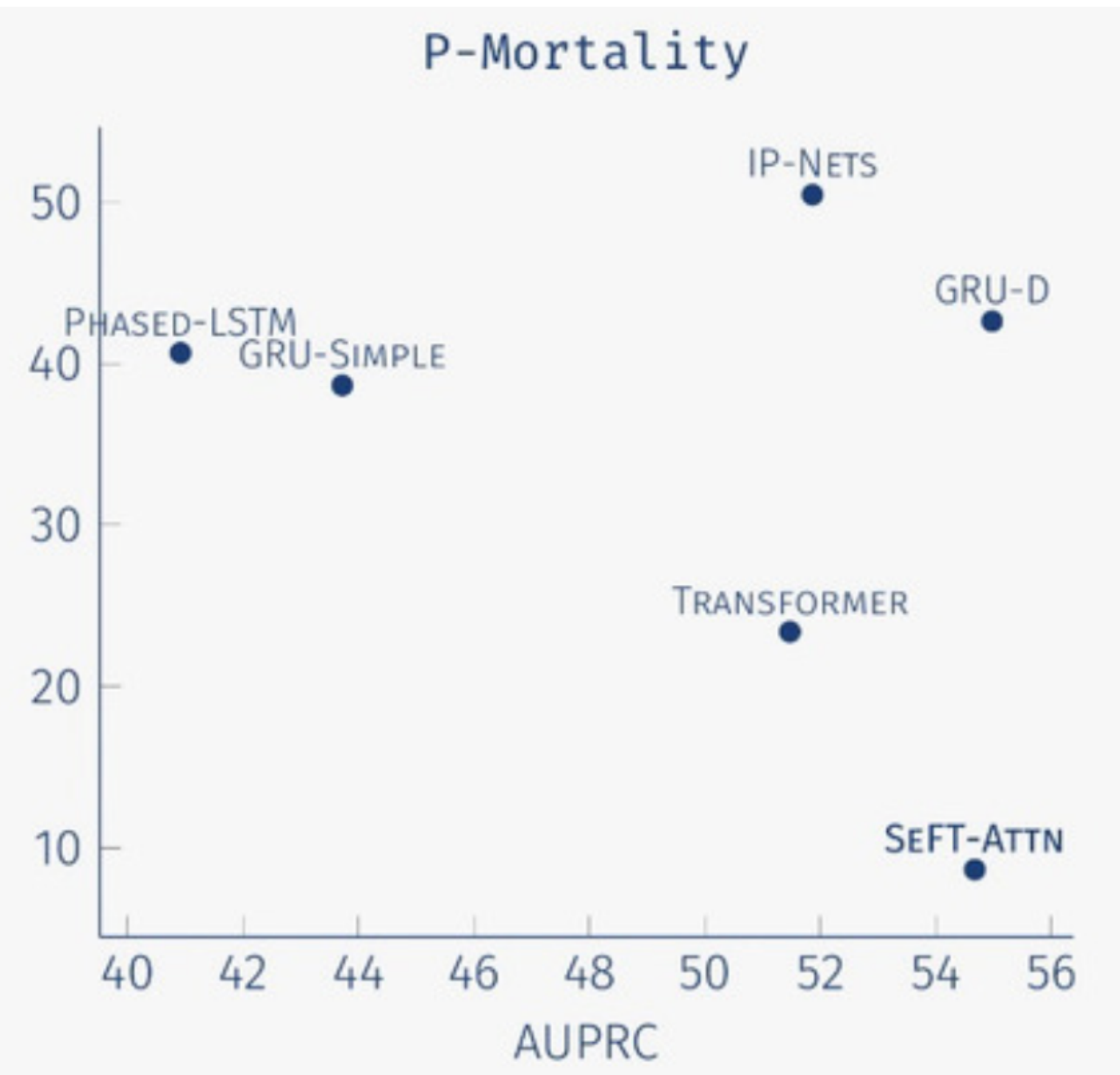
Dataset	Model	Accuracy	AUPRC	AUROC	s/epoch
M3M	GRU-D	77.0 $\pm$ 1.5	<b>52.0 <math>\pm</math> 0.8</b>	<b>85.7 <math>\pm</math> 0.2</b>	133 $\pm$ 8
	GRU-SIMPLE	78.1 $\pm$ 1.3	43.6 $\pm$ 0.4	82.8 $\pm$ 0.0	140 $\pm$ 7
	IP-NETS	78.3 $\pm$ 0.7	48.3 $\pm$ 0.4	83.2 $\pm$ 0.5	81.2 $\pm$ 8.5
	PHASED-LSTM	73.8 $\pm$ 3.3	37.1 $\pm$ 0.5	80.3 $\pm$ 0.4	166 $\pm$ 7
	TRANSFORMER	77.4 $\pm$ 5.6	42.6 $\pm$ 1.0	82.1 $\pm$ 0.3	20.1 $\pm$ 0.1
	LATENT-ODE <sup>†</sup>	72.8 $\pm$ 1.7	39.5 $\pm$ 0.5	80.9 $\pm$ 0.2	4622
	SEFT-ATTN	<b>79.0 <math>\pm</math> 2.2</b>	46.3 $\pm$ 0.5	83.9 $\pm$ 0.4	<b>14.5 <math>\pm</math> 0.5</b>
P12	GRU-D	80.0 $\pm$ 2.9	53.7 $\pm$ 0.9	<b>86.3 <math>\pm</math> 0.3</b>	8.67 $\pm$ 0.49
	GRU-SIMPLE	82.2 $\pm$ 0.2	42.2 $\pm$ 0.6	80.8 $\pm$ 1.1	30.0 $\pm$ 2.5
	IP-NETS	79.4 $\pm$ 0.3	51.0 $\pm$ 0.6	86.0 $\pm$ 0.2	25.3 $\pm$ 1.8
	PHASED-LSTM	76.8 $\pm$ 5.2	38.7 $\pm$ 1.5	79.0 $\pm$ 1.0	44.6 $\pm$ 2.3
	TRANSFORMER	<b>83.7 <math>\pm</math> 3.5</b>	<b>52.8 <math>\pm</math> 2.2</b>	<b>86.3 <math>\pm</math> 0.8</b>	<b>6.06 <math>\pm</math> 0.06</b>
	LATENT-ODE <sup>†</sup>	76.0 $\pm$ 0.1	50.7 $\pm$ 1.7	85.7 $\pm$ 0.6	3500
	SEFT-ATTN	75.3 $\pm$ 3.5	52.4 $\pm$ 1.1	85.1 $\pm$ 0.4	7.62 $\pm$ 0.10

# SeFT Result: Performance vs. Runtime

MIMIC-III Mortality Prediction Task



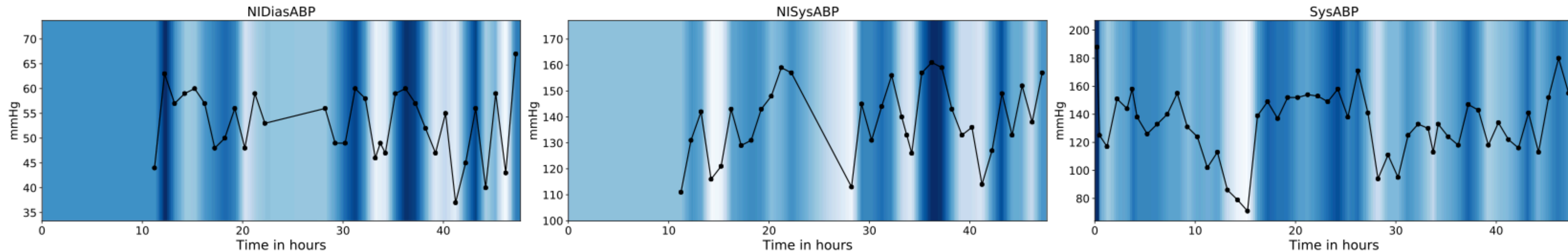
Physionet 2012 Mortality Prediction Task





# SeFT's Outputs are Explainable

---

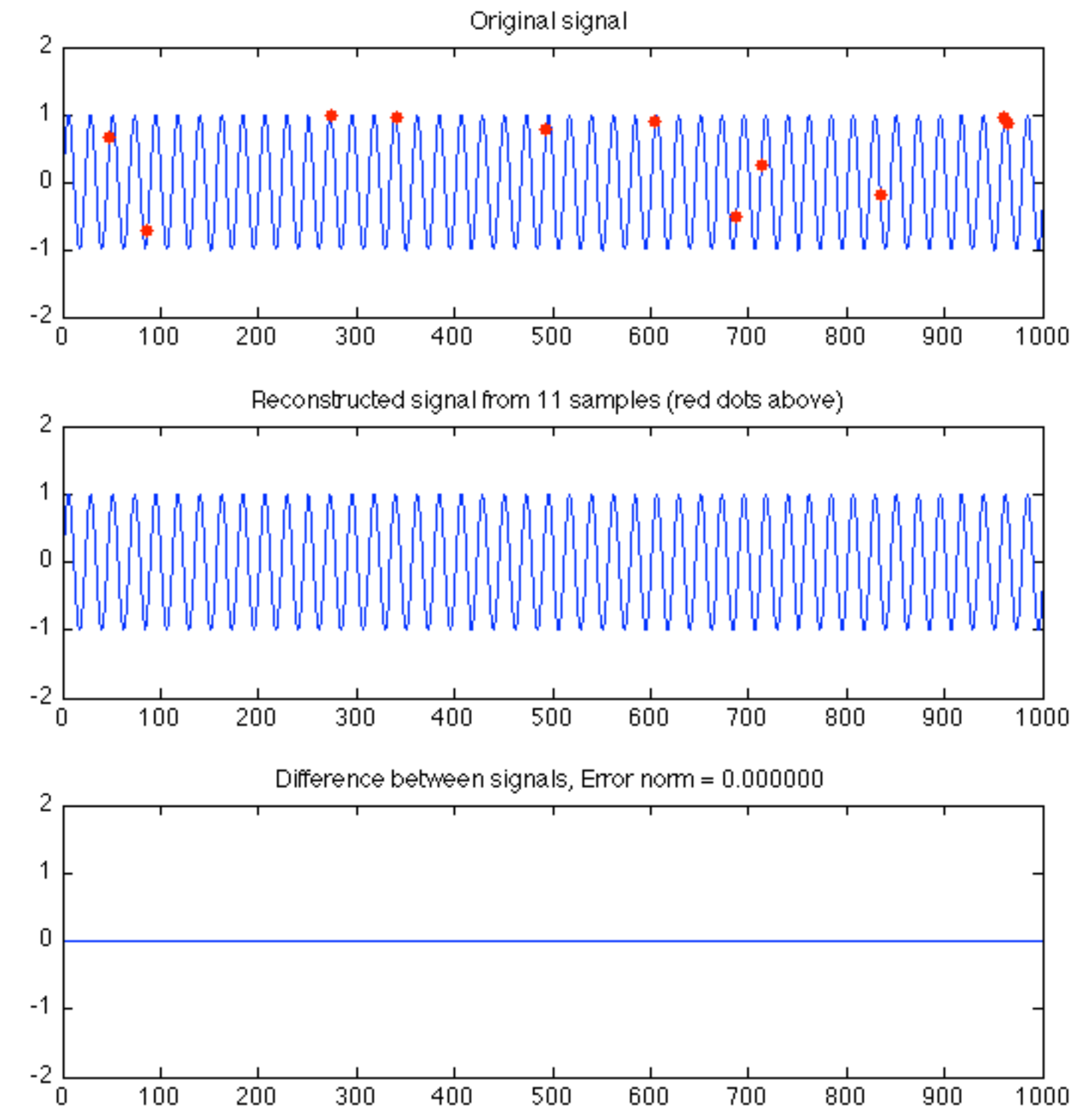
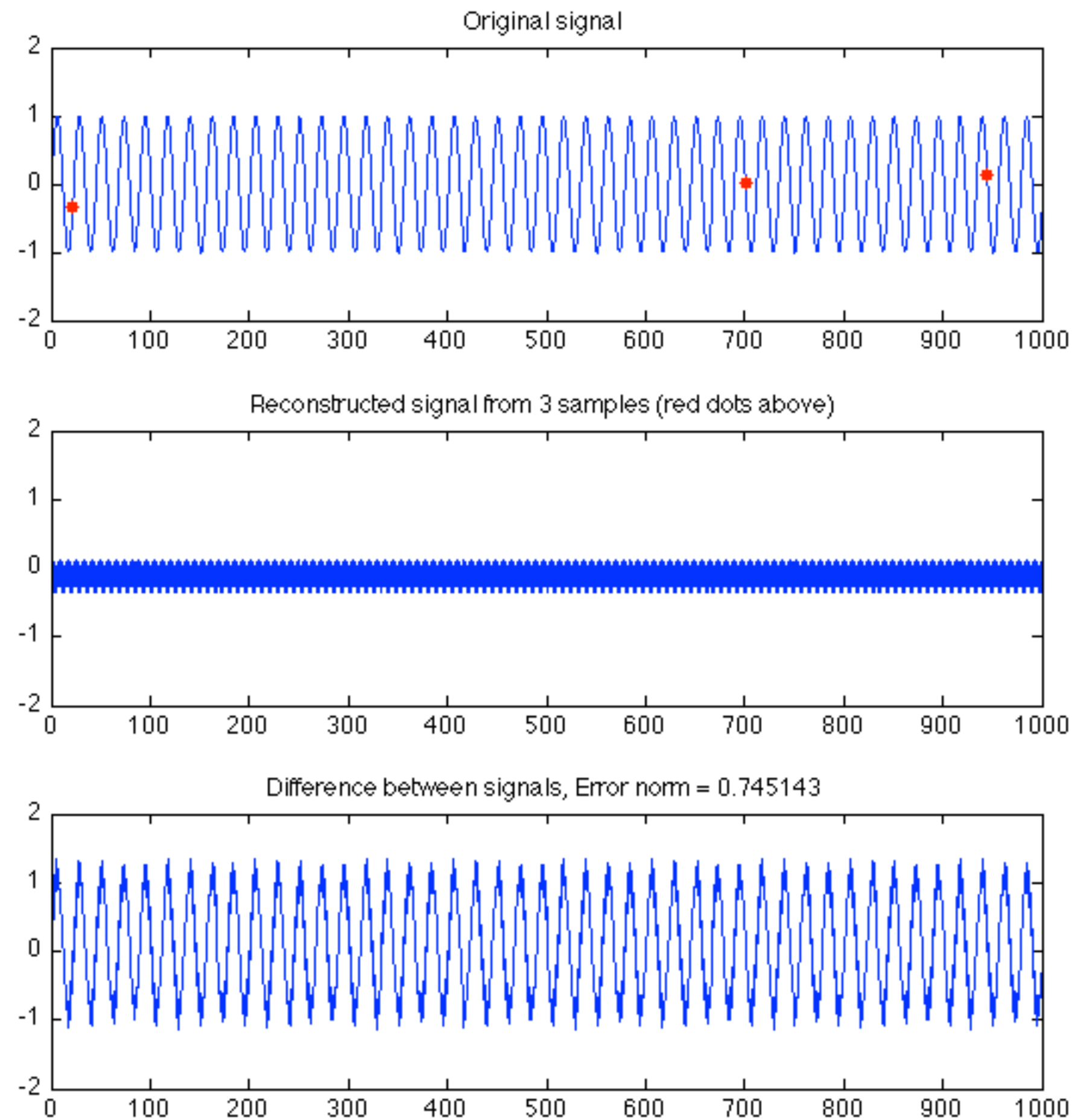


Allows a per-observation quantification of importance

# Combining Compressed Sensing and Machine Learning



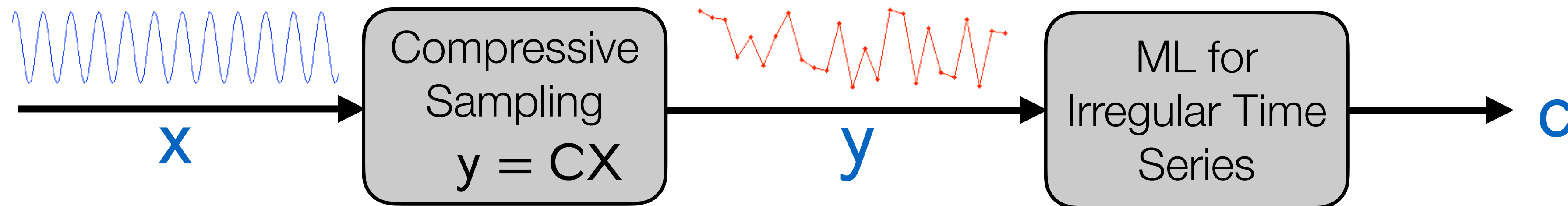
# Recall: Compressed Sensing



# Combining Compressive Sensing and Machine Learning

---

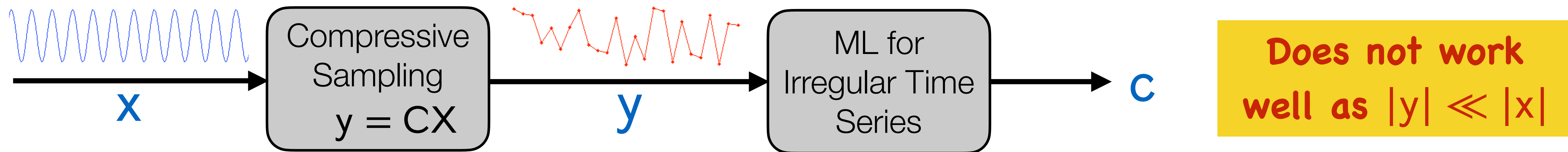
## Option 1



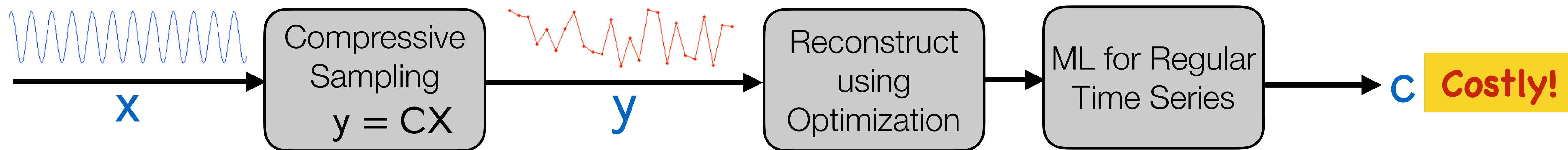
**Does not work  
well as  $|y| \ll |x|$**

# Combining Compressive Sensing and Machine Learning

## Option 1



## Option 2



$$\operatorname{argmin}_s \|s\|_1, \text{ s.t. } \|C\Psi s - y\|_2 < \epsilon$$



**Could we Combine  
Compressive Sensing and  
Deep Learning?**

# Compressed Learning with Deep Neural Networks

---

# Compressed Learning with Deep Neural Networks

---

- Compressed Learning: direct inference from compressive measurements is feasible with high classification accuracies
  - ▶ Calderbank, Jafarpour, and Schapire. "Compressed learning: Universal sparse dimensionality reduction and learning in the measurement domain." preprint (2009).
  - ▶ Proved that under certain conditions the performance of a linear SVM classifier operating in the compressed sensing domain  $y = Cx$  is almost equivalent to the performance of the best linear threshold classifier operating in the signal domain  $x$

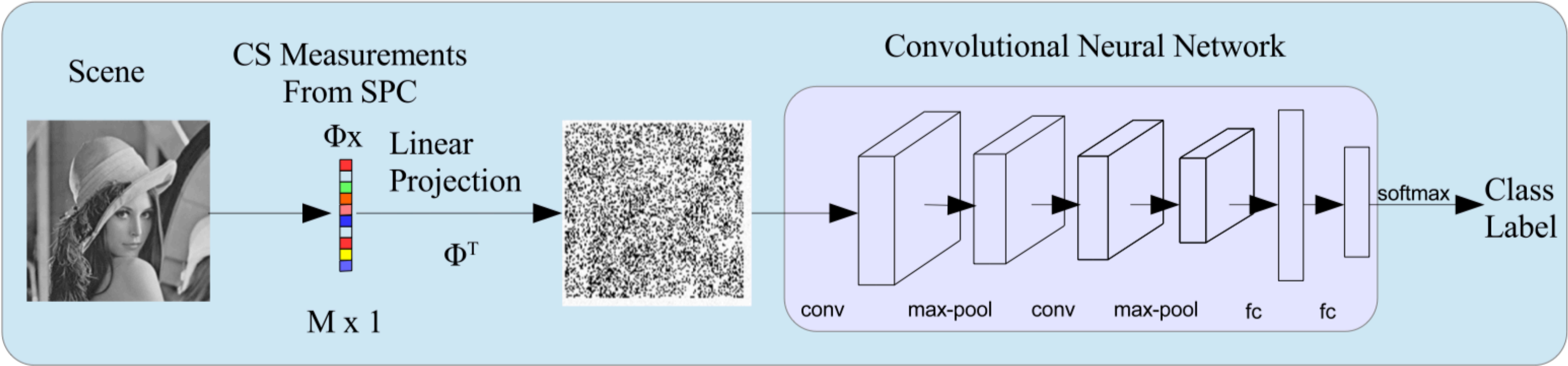


# Compressed Learning with Deep Neural Networks

---

- Compressed Learning: direct inference from compressive measurements is feasible with high classification accuracies
  - ▶ Calderbank, Jafarpour, and Schapire. "Compressed learning: Universal sparse dimensionality reduction and learning in the measurement domain." preprint (2009).
  - ▶ Proved that under certain conditions the performance of a linear SVM classifier operating in the compressed sensing domain  $y = Cx$  is almost equivalent to the performance of the best linear threshold classifier operating in the signal domain  $x$
- Subsequent work combined CL with DNNs: projected measurement vector  $z = C^T y$  (same shape as  $x$ ) as the input to a convolutional neural network
  - ▶ Lohit, Kulkarni, and Turaga. "Direct inference on compressive measurements using convolutional neural networks." IEEE Intl Conf on Image Proc. (ICIP), pp. 1913-1917.
  - ▶ Showed pretty good results on MNIST and ImageNet classification by training on projected measurement  $z = C^T y$  instead of on the original signal  $x$

# Direct Inference on Compressive Measurements using CNN



Note: specific CNN architecture is for illustration only - specific architecture will depend on application.

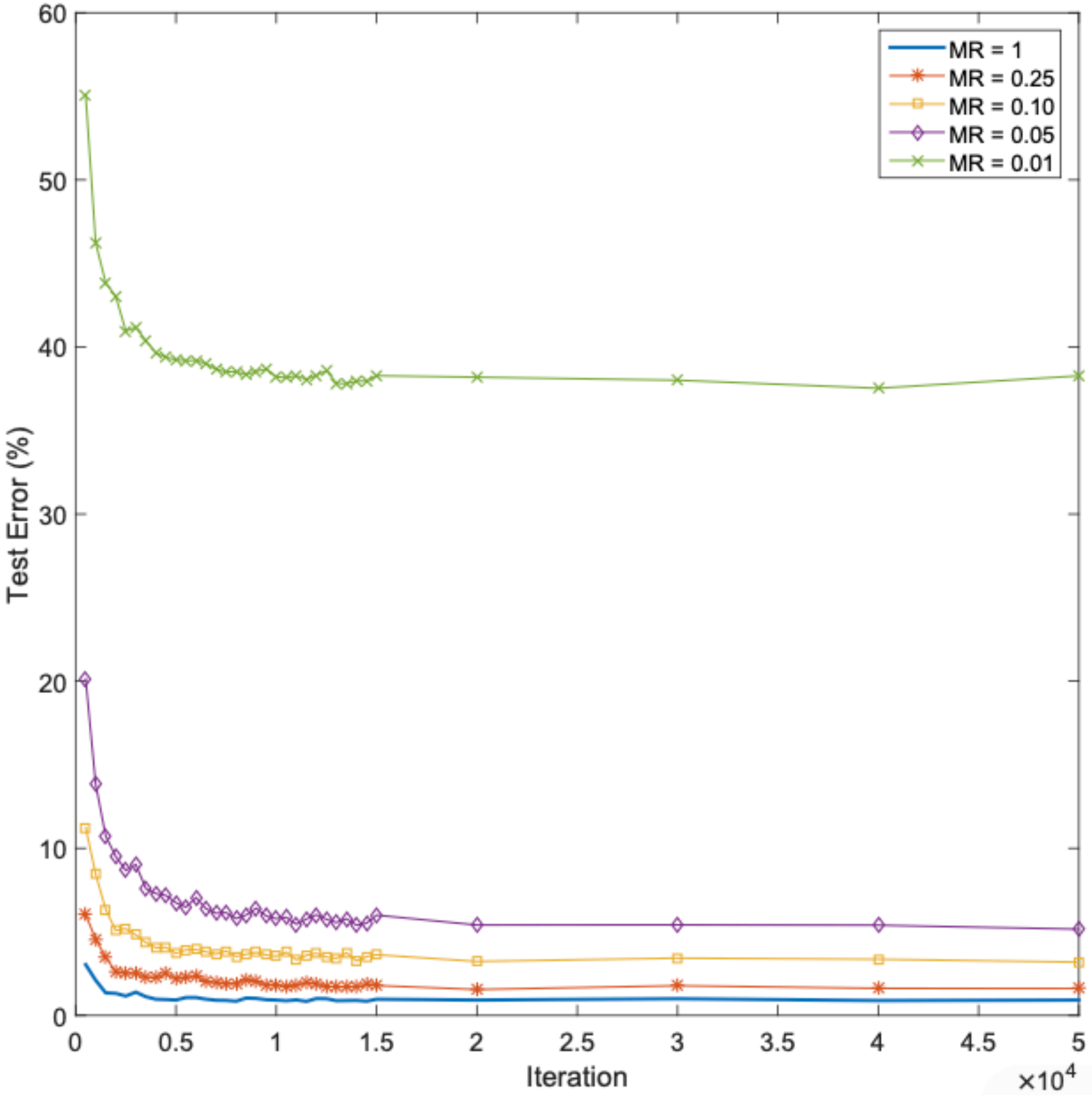
## Architecture

Measurement Rate (MR)	Number of Measurements	Test Error	
		Smashed Filters [4]	Our Method
1 (Oracle)	784	13.86%	0.89%
0.25	196	27.42%	1.63%
0.10	78	43.55%	2.99%
0.05	39	53.21%	5.18%
0.01	8	63.03%	41.06%

## MNIST Classification

Measurement Rate	No. of Measurements	Accuracy
1 (Oracle)	65536	56.88%
0.25	16384	39.22%
0.10	6554	29.84%

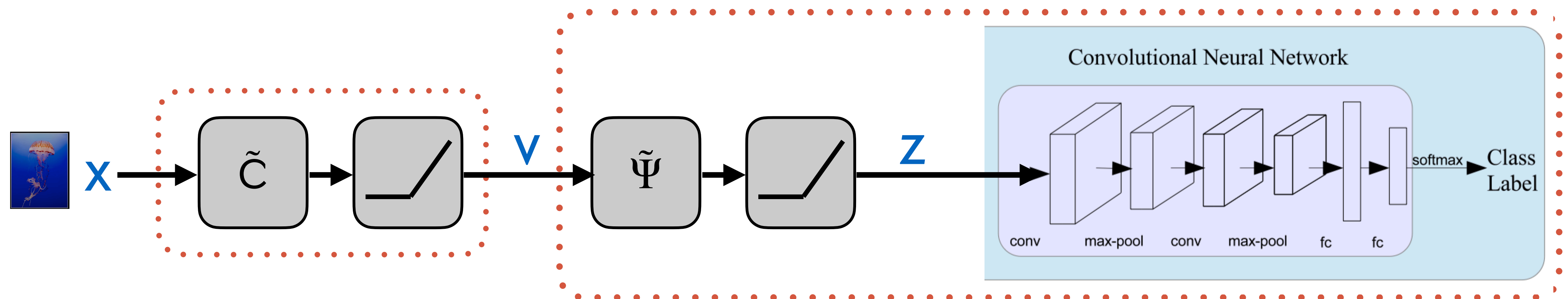
## ImageNet Classification



## Convergence of Test Error

# End-to-end Deep Learning Solution for Compressed Learning

- Idea: jointly optimize the sensing matrix  $C$  and the inference operator (i.e. the CNN)
  - ▶ Adler, Amir, Michael Elad, and Michael Zibulevsky. "Compressed learning: A deep neural network approach." arXiv preprint arXiv:1610.09615 (2016).
- Approach
  - ▶ The first layer learns and performs the sensing matrix  $C$
  - ▶ The subsequent layers (a fully-connected layer followed by a CNN) perform the non-linear inference stage
  - ▶ The second fully-connected layer performs operation similar to  $z = C^T y$  but a different matrix  $\tilde{C}$  is learnt
  - ▶ The first and second layers are followed by ReLU
  - ▶ The two components of end-to-end CL detached after training





# End-to-end Deep Learning Solution for Compressed Learning

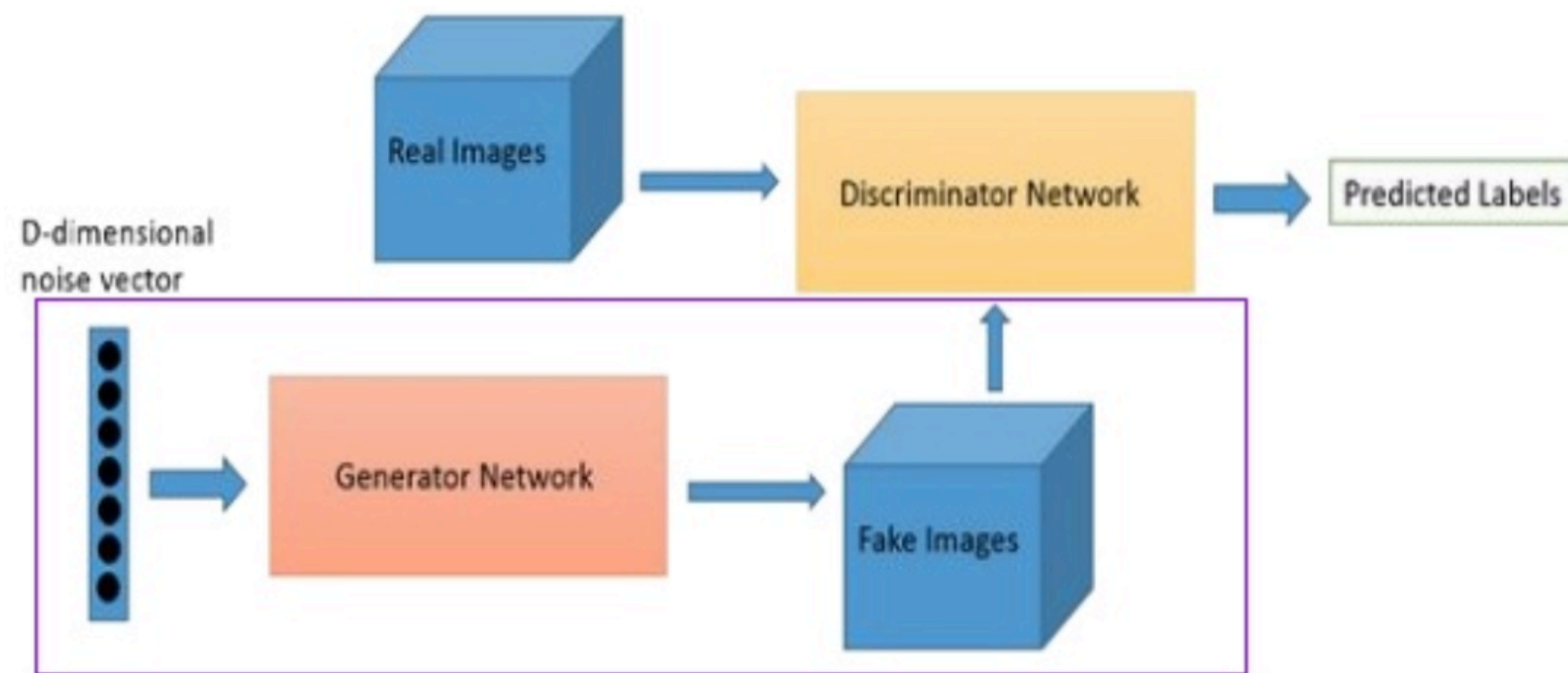
- Idea: jointly optimize the sensing matrix  $C$  and the inference operator (i.e. the CNN)
  - ▶ Adler, Amir, Michael Elad, and Michael Zibulevsky. "Compressed learning: A deep neural network approach." arXiv preprint arXiv:1610.09615 (2016).
- Approach
  - ▶ The first layer learns and performs the sensing matrix  $C$
  - ▶ The subsequent layers (a fully-connected layer followed by a CNN) perform the non-linear inference stage
  - ▶ The second fully-connected layer performs operation similar to  $z = C^T y$  but a different matrix  $\tilde{C}$  is learnt
  - ▶ The first and second layers are followed by ReLU
  - ▶ The two components of end-to-end CL detached after training

Sensing Rate	No. of Measurements	Smashed Filters [12]	Random Sensing + CNN [4]	Proposed
0.25	196	27.42%	1.63%	<b>1.56%</b>
0.1	78	43.55%	2.99%	<b>1.91%</b>
0.05	39	53.21%	5.18%	<b>2.86%</b>
0.01	8	63.03%	41.06%	<b>6.46%</b>

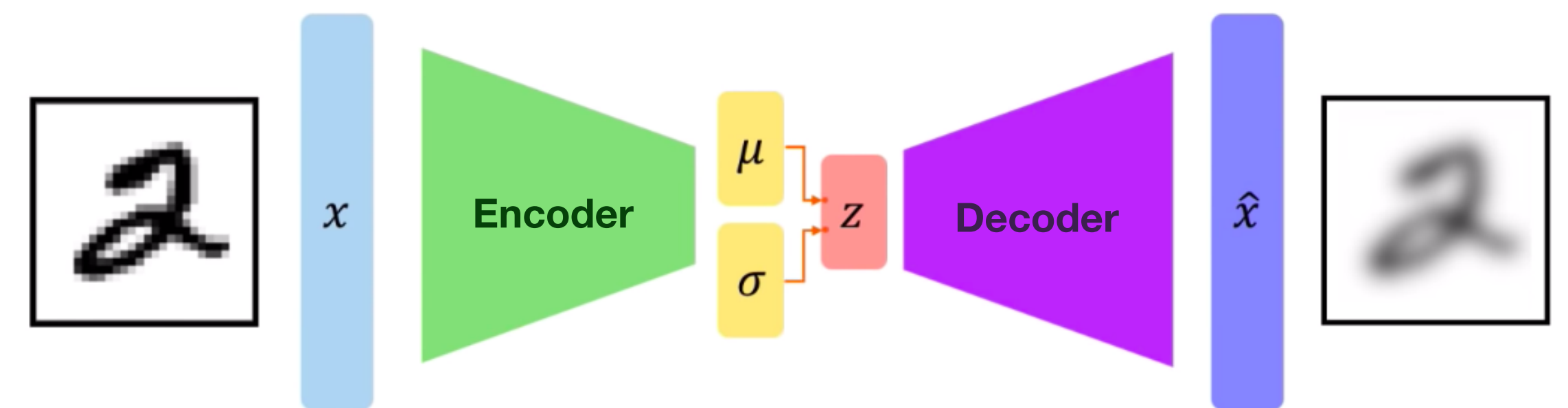
Classification Error (%) for the MNIST handwritten digits dataset vs. sensing rate  $R = M/N$  (averaged over 10,000 test images)

# Compressed Sensing Using Generative Models

- CS compresses by taking random linear projections (measurement matrix) of the original signal, and reconstructs by exploiting sparsity present in “natural” signals
- Instead of relying on sparsity, one can use *structure* from a generative model.
  - GANs and VAEs



**Generative Adversarial Networks**



**Variational Auto Encoder**



# Compressed Sensing Using Generative Models

---

- CS compresses by taking random linear projections (measurement matrix) of the original signal, and reconstructs by exploiting sparsity present in “natural” signals
- Instead of relying on sparsity, one can use *structure* from a generative model.
  - ▶ GANs and VAEs
  - ▶ A generative model is given by a deterministic function  $G : \mathbb{R}^k \rightarrow \mathbb{R}^n$ , and a distribution  $P_Z$  over  $z \in \mathbb{R}^k$ .
  - ▶ To generate a sample from the generator, we draw  $z \sim P_Z$  and the sample then is  $G(z)$
- **Approach:** find a vector in representation space s.t. the corresponding vector in the sample space matches the observed measurements, i.e. optimize  $loss(z) = ||AG(z) - y||_2^2$  (highly non-convex, approximated using gradient descent)
  - ▶ If the optimization procedure gives  $\hat{z}$ , then reconstruct by computing  $\hat{x} = G(\hat{z})$

# Compressed Sensing Using Generative Models

- CS compresses by taking random linear projections (measurement matrix) of the original signal, and reconstructs by exploiting sparsity present in “natural” signals

- Instead

- ▶ GAN

- ▶ A generator

- over

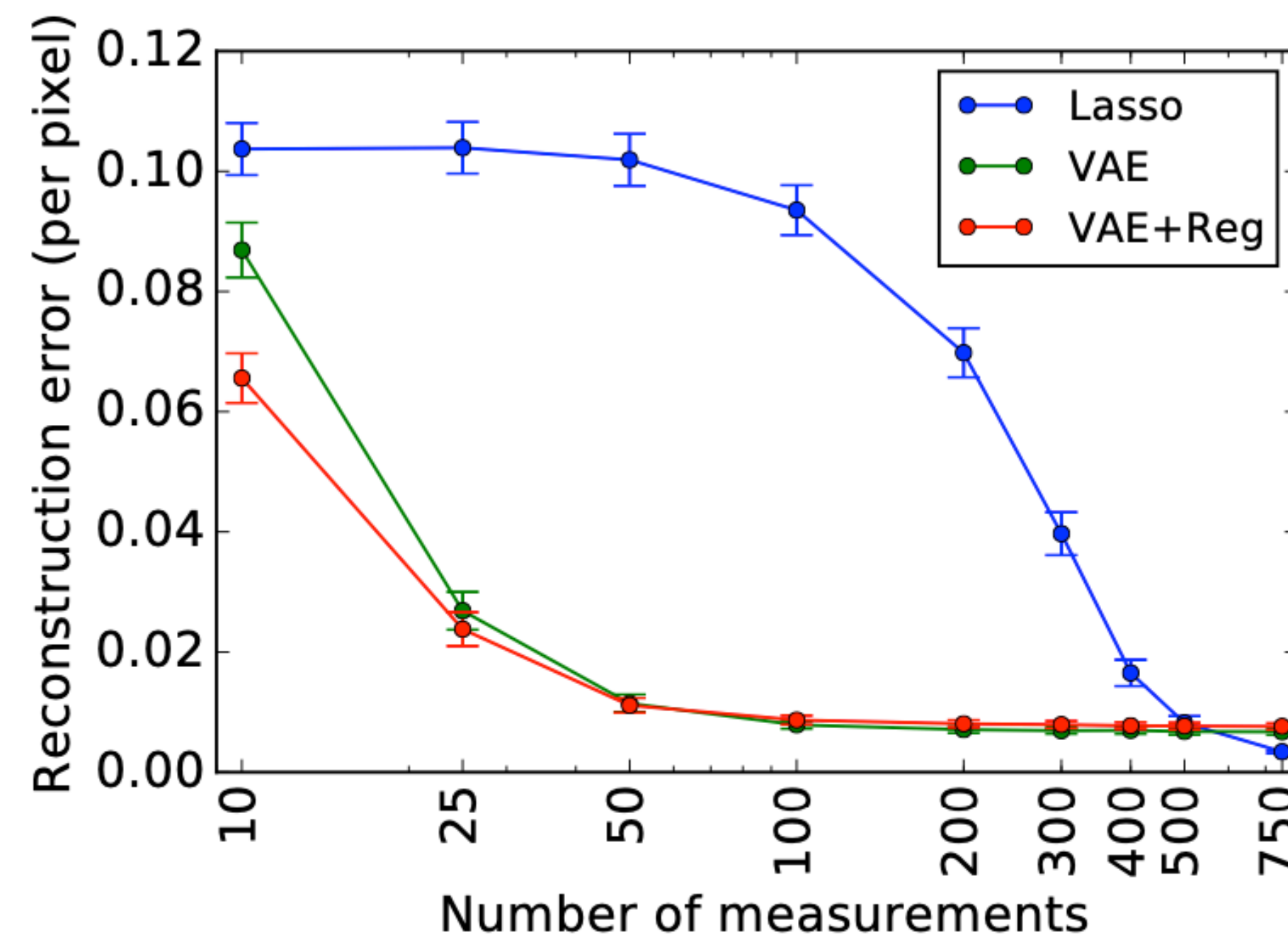
- ▶ To generate

- Application

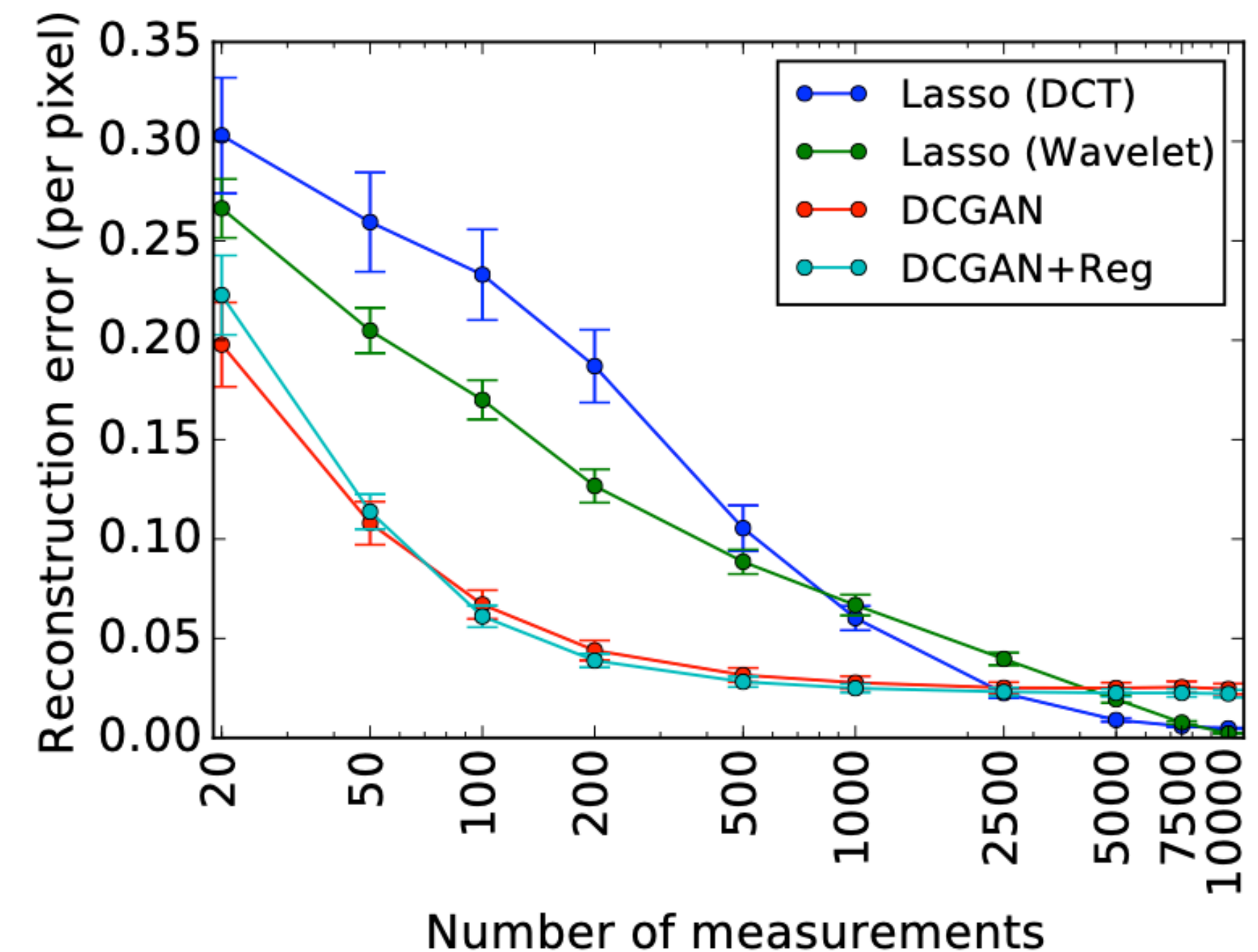
- signal

- loss

- ▶ If the optimization procedure gives  $\hat{z}$ , then reconstruct by computing  $\hat{x} = G(\hat{z})$



(a) Results on MNIST



(b) Results on celebA

# References

---

- Zaheer, Manzil, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Ruslan Salakhutdinov, and Alexander Smola. "Deep sets." arXiv preprint arXiv:1703.06114 (2017).
- Chen, Ricky TQ, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. "Neural ordinary differential equations." arXiv preprint arXiv:1806.07366 (2018).
- Shukla, Satya Narayan, and Benjamin M. Marlin. "Interpolation-prediction networks for irregularly sampled time series." arXiv preprint arXiv:1909.07782 (2019).
- Rubanova, Yulia, Ricky TQ Chen, and David Duvenaud. "Latent ODEs for irregularly-sampled time series." arXiv preprint arXiv:1907.03907 (2019).
- Horn, Max, Michael Moor, Christian Bock, Bastian Rieck, and Karsten Borgwardt. "Set functions for time series." In International Conference on Machine Learning, pp. 4353-4363. PMLR, 2020.
- Shukla, Satya Narayan, and Benjamin M. Marlin. "A Survey on Principles, Models and Methods for Learning from Irregularly Sampled Time Series: From Discretization to Attention and Invariance." arXiv preprint arXiv:2012.00168 (2020).
- Shukla, Satya Narayan, and Benjamin M. Marlin. "Multi-Time Attention Networks for Irregularly Sampled Time Series." arXiv preprint arXiv:2101.10318 (2021).

The End